# Biostatistics for Predicting Health Disparities in Infectious Disease Outcomes, Using Real-world Evidence and Public Health Intervention Data

DAVID OCHE IDOKO[1], OKOROJI EMMANUEL MBACHU[2], IDAYAT NINILOLA OLOLADE BABALOLA[3], ERONDU OKECHUKWU FELIX[4], OLUWAYEMISI DADA-ABIDAKUN[5], YEWANDE ADEYEYE[6]

[1]*Department of Fisheries and Aquaculture, J.S Tarkaa University, Makurdi, Nigeria.*

[2]*Department of Obstetrics and Gynecology, David Umahi Federal University Teaching Hospital, Uburu, Ebonyi State, Nigeria.*

[3]*Supported Living Services, Time 4 U Ltd, Chatham, UK.*

[4]*Department of Radiography and Radiation Sciences, Gregory University, Uturu, Abia State, Nigeria.*

[5]*Federal Teaching Hospital, Ado Ekiti, Ekiti, Nigeria.*

[6]*Day Case Surgery Department, Warrington and Halton Hospital, Warrington City, United Kindom.*

*Abstract- This review explores emerging biostatistical methods, the integration of machine learning (ML) and advanced analytics, and the role of big data and artificial intelligence (AI) in addressing health disparities in public health. It highlights the growing importance of Bayesian models and ML algorithms for predicting infectious disease outcomes and stratifying populations by social determinants of health. The review accentuates the potential of AI in precision public health, with applications ranging from real-time disease surveillance to the development of personalized interventions. However, it also emphasizes the ethical challenges and biases associated with AI and ML, particularly in marginalized populations. Future research recommendations focus on developing ethical frameworks, improving the representativeness of training data, and optimizing the use of real-world evidence (RWE) in public health. By combining traditional biostatistical approaches with modern AI-driven tools, this review outlines a path toward more accurate and equitable health outcome predictions, ultimately contributing to the reduction of health disparities on a global scale.*

*Indexed Terms- Biostatistical methods, Machine learning, Health disparities, Artificial intelligence, Real-world evidence*

## I. INTRODUCTION

### 1.1 Overview of Health Disparities in Infectious Disease Outcomes

Health disparities in infectious disease outcomes are well-documented and persist across various populations, often due to differences in socioeconomic factors, access to healthcare, and biological susceptibilities (Marmot, 2005). Infectious diseases such as tuberculosis, HIV/AIDS, and malaria disproportionately affect marginalized communities, both in high-income countries and low- and middle-income regions (WHO, 2021). For instance, in the United States, Black and Hispanic populations are more likely to experience higher rates of HIV infection compared to White populations, with an estimated rate of 41.3 per 100,000 among Black individuals, compared to 5.0 per 100,000 for White individuals (CDC, 2022). This disparity is often linked to structural inequities, including poverty, lack of access to preventive healthcare services, and stigmatization (Kawachi et al., 2002). The COVID-19 pandemic further illustrated these inequalities, as minority populations in the United States and globally experienced higher rates of infection, hospitalization, and mortality (Bambra et al., 2020).
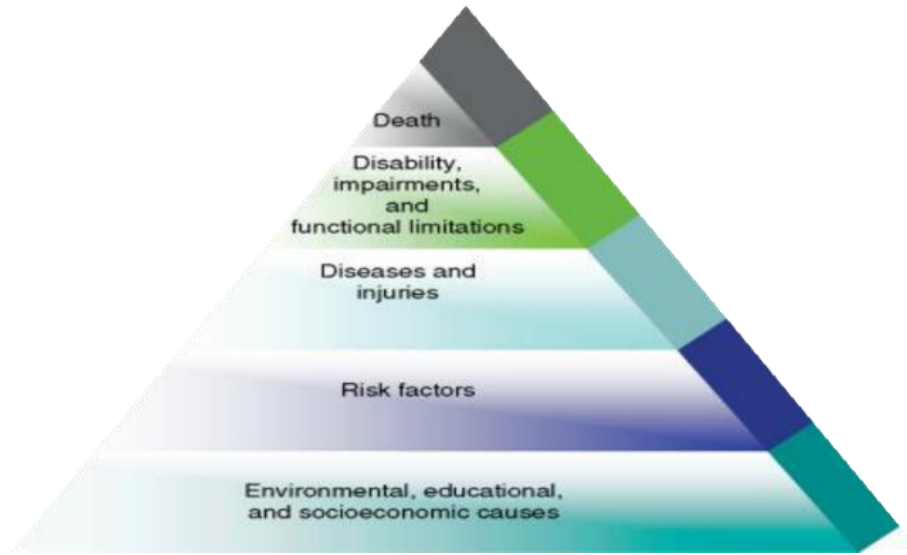
Figure 1: Overview of the burden of disease framework. (Bhutta, 2008)

The figure above illustrates the progression and interconnectedness of health issues and their primary causes. It implies a hierarchical relationship, where each level is influenced by the ones below it. It also suggests that addressing issues at the lower levels could have a cascading positive effect on the levels above. The color gradient from teal at the bottom to gray at the top adds visual clarity to the progression of health issues from underlying causes to ultimate outcomes.

Biological factors also play a role in these disparities. Genetic predispositions, immune response variations, and comorbidities such as diabetes and cardiovascular diseases exacerbate the risk of poor outcomes in certain populations (Dowd et al., 2009). For example, research has shown that individuals with compromised immune systems or underlying chronic conditions are more likely to experience severe complications from infectious diseases like influenza and COVID-19 (Sattar et al., 2020). Environmental factors, including housing conditions, pollution exposure, and occupational hazards, further compound the vulnerability of disadvantaged groups to infectious diseases (Phelan et al., 2010). These intertwined factors necessitate a multifaceted approach to addressing health disparities, including improved access to healthcare, targeted public health interventions, and comprehensive biostatistical analyses to predict and mitigate the risks faced by vulnerable populations.

Real-world evidence suggests that public health interventions tailored to address specific social determinants of health can help reduce disparities in infectious disease outcomes. For example, increasing access to vaccines, improving sanitation, and implementing educational campaigns in underserved communities have been shown to lower infection rates and improve health outcomes (Farmer et al., 2006). However, addressing these disparities requires a sustained effort to integrate biostatistics, public health strategies, and equitable healthcare delivery systems. By identifying the key drivers of these inequalities, policymakers and healthcare providers can better design interventions that specifically target the populations most at risk, thus improving overall public health outcomes.

1.2 Importance of biostatistics in addressing these disparities
Biostatistics plays a pivotal role in addressing health disparities in infectious disease outcomes by providing the tools necessary to analyze complex datasets and uncover patterns that may not be immediately visible. Through statistical modeling, biostatistics enables researchers to identify correlations between demographic factors—such as race, socioeconomic status, and geographic location—and infectious

disease incidence and outcomes (Diez Roux, 2012). By analyzing large- scale data, such as real-world evidence from electronic health records and national health surveys, biostatistical methods help quantify the extent of disparities and assess the effectiveness of interventions targeted at vulnerable populations (Rosella et al., 2018). For instance, regression models have been widely used to control for confounding factors and estimate the relative risk of infection or poor health outcomes, thereby isolating the specific contribution of social determinants of health to disease disparities (Vandenbroucke et al., 2007).

Furthermore, biostatistics facilitates the measurement of the impact of public health interventions on reducing disparities. Randomized controlled trials (RCTs) and observational studies rely heavily on biostatistical techniques to evaluate the efficacy of interventions, such as vaccination campaigns or community-based health programs, in different population subgroups (Pocock, 2013). Through the application of survival analysis, for example, researchers can track long-term outcomes of patients from disadvantaged communities and monitor the effectiveness of preventive

measures over time (Lai et al., 2021). Additionally, biostatistical methods are essential in addressing the biases that often arise in real-world datasets, such as missing data or selection bias, which can distort the findings if not properly accounted for (Rubin, 2004).

The ability of biostatistics to synthesize data from diverse sources and correct for such biases ensures that the results are robust and generalizable across populations. This allows policymakers and public health officials to allocate resources more effectively and design evidence-based interventions that specifically target the most vulnerable groups. Moreover, biostatistical methods enable the exploration of interaction effects, such as how the combination of low socioeconomic status and inadequate healthcare access exacerbates disease risks, thereby providing a more nuanced understanding of the multifactorial nature of health disparities (Subramanian and Kawachi, 2004). Overall, biostatistics is an indispensable tool in the fight against health inequities, offering data-driven

insights that can be translated into actionable public health strategies.

1.3     The role of real-world evidence and public health intervention data

Real-world evidence (RWE) has become a crucial component in understanding and addressing health disparities in infectious disease outcomes. Derived from sources such as electronic health records, insurance claims, and health surveys, RWE provides a more comprehensive view of how diseases affect different populations in naturalistic settings (Makady et al., 2017). Unlike data obtained from controlled clinical trials, which often exclude vulnerable populations or lack diversity, real-world evidence reflects the actual experiences of patients, including those from disadvantaged communities (Corrigan-Curay et al., 2018). This data is invaluable in assessing health disparities, as it enables researchers to identify patterns of disease progression, treatment responses, and the effectiveness of public health interventions across diverse socioeconomic and racial groups (Sherman et al., 2016). For example, in a large-scale study of influenza vaccination, RWE revealed that minority populations had lower vaccination rates, contributing to higher hospitalization rates among these groups during the flu season (Lu et al., 2014).

Public health intervention data complements real-world evidence by providing insight into the effectiveness of strategies aimed at reducing health disparities. This data is typically gathered from community-level programs, national health campaigns, and localized interventions aimed at improving disease outcomes (Galea et al., 2019). By combining RWE with public health intervention data, biostatisticians can measure the long-term impact of interventions, such as vaccination drives, sanitation improvements, or health education campaigns, on population health outcomes (Pelat et al., 2014). For example, RWE from public health interventions targeting tuberculosis in low-income areas has demonstrated significant reductions in disease prevalence when resources such as early diagnosis, contact tracing, and treatment adherence programs are implemented. Such data-driven insights are critical in optimizing public health strategies, allowing interventions to be tailored to the specific needs of vulnerable populations and adjusted in real- time based on their effectiveness.
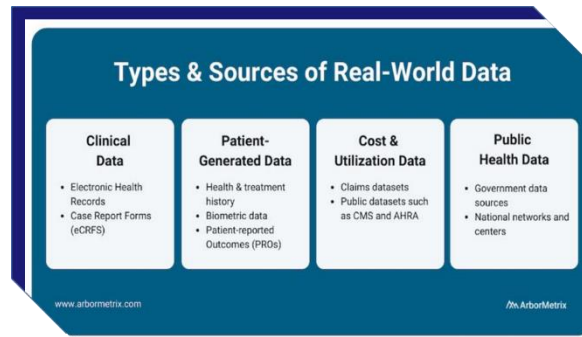
Figure 2: Types & Sources of Real-World Data (Arbor, 2021)

This figure effectively illustrates the diverse range of data sources available for real-world evidence in healthcare, highlighting the multipart nature of health-related information. This type of data is crucial for comprehensive healthcare research, policy-making, and improving patient outcomes.

The role of real-world evidence and public health intervention data extends beyond merely identifying disparities. These data sources are integral to policy formulation, as they provide empirical support for the allocation of resources and the design of equitable healthcare solutions (Harron et al., 2017). With accurate and robust real-world evidence, public health officials can develop interventions that are both effective and cost-efficient, ensuring that high-risk populations receive the care and prevention strategies they require. Moreover, the use of biostatistics in synthesizing data from these diverse sources allows for the identification of subtle yet critical factors—such as cultural barriers or healthcare access inequalities—that may otherwise be overlooked (Idoko et al., 2024). By focusing on real-world evidence and public health intervention data, healthcare systems can become more responsive and adaptive, ultimately contributing to the reduction of health disparities in infectious disease outcomes.

### 1.4 Objectives of the review paper

The primary objective of this review paper is to explore the role of biostatistics in predicting health disparities in infectious disease outcomes, particularly through the analysis of real-world evidence and public health intervention data. By synthesizing existing literature and empirical studies, this review aims to provide a comprehensive understanding of how

biostatistical methods can help identify and address the social determinants of health that contribute to these disparities (Diez Roux, 2012). A secondary objective is to highlight the importance of data-driven public health interventions, demonstrating how real-world evidence can be leveraged to tailor strategies that reduce infection rates and improve outcomes in marginalized populations (Rosella et al., 2018). This review will focus on the use of advanced biostatistical techniques, including regression models and survival analysis, to evaluate the long-term impact of these interventions.

Another key objective is to examine how biostatistics can contribute to equitable healthcare by ensuring that public health resources are efficiently allocated to those most at risk (Idoko et al., 2024). Through a detailed analysis of public health intervention data, this paper seeks to identify the specific factors—such as vaccination coverage, healthcare access, and socioeconomic conditions—that exacerbate disparities in infectious disease outcomes (Bambra et al., 2020). Ultimately, the paper will argue that integrating biostatistics into public health planning is essential for mitigating the disproportionate burden of infectious diseases on disadvantaged communities.

In fulfilling these objectives, this review will also address gaps in the current literature by assessing the methodological challenges involved in using real-world evidence and public health data. It aims to offer recommendations for future research and public health policy, encouraging the continued development of biostatistical models that are responsive to the unique health needs of vulnerable populations (Pocock, 2013). This review will thus contribute to the broader

effort of reducing health disparities and improving disease outcomes through the rigorous application of biostatistics in public health.

## II. BIOSTATISTICAL METHODS IN PREDICTING HEALTH DISPARITIES

### 2.1 Key biostatistical tools used in health disparity prediction

In the prediction of health disparities, several biostatistical tools are fundamental in analyzing large and complex datasets, enabling researchers to discern patterns and trends across different population groups. One of the most widely used techniques is logistic regression, which is crucial for modeling binary outcomes, such as the presence or absence of a disease, and identifying the influence of demographic and socioeconomic variables on these outcomes (Hosmer et al., 2013). Logistic regression allows for the estimation of odds ratios, which quantify the association between specific risk factors and disease outcomes, making it highly applicable for assessing the likelihood of adverse health outcomes in disadvantaged populations (Diez Roux, 2012). By controlling for confounding variables, this tool can isolate the effect of race, income, or geographic location on health disparities in infectious diseases such as tuberculosis and HIV/AIDS (Koh et al., 2012).
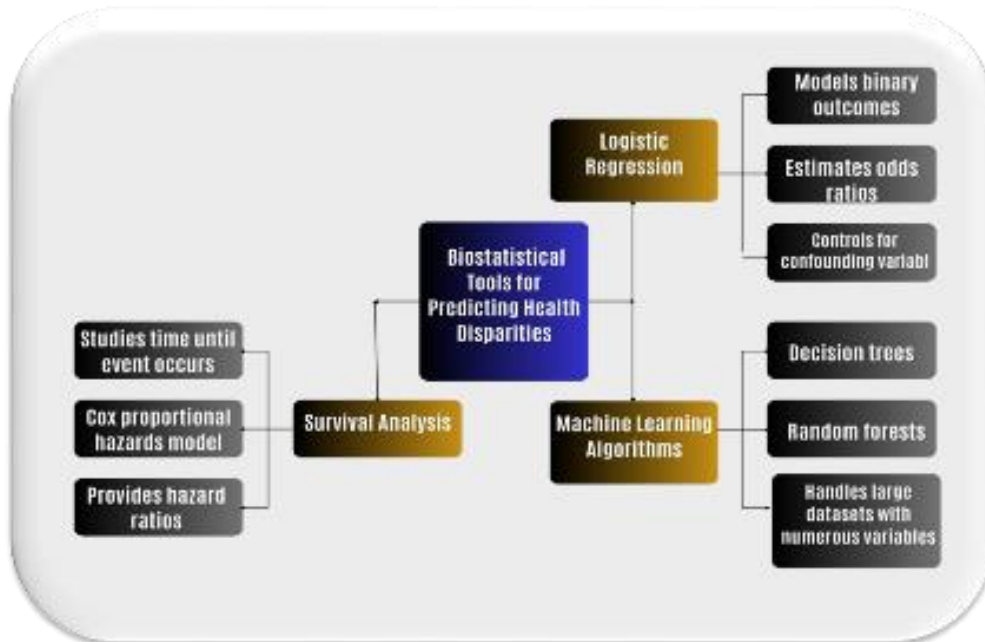


Figure 3: Biostatistical Tools for Predicting Health Disparities

This visual representation above shows the relationships between different biostatistical tools and their applications in health disparity research. It captures the complex nature of the statistical approaches used in this field.

Survival analysis is another key biostatistical method, particularly useful for studying the time until an event occurs, such as death or disease progression (Collett, 2015). In health disparity research, survival models, like the Cox proportional hazards model, are employed to assess how long different population groups remain disease-free or survive after diagnosis, accounting for censored data and covariates (Hosmer et al., 2008). This is particularly valuable in examining the long-term impact of social determinants of health, such as access to healthcare or education, on infectious disease outcomes. For instance, survival analysis has been used to demonstrate that lower socioeconomic status is associated with reduced survival rates in diseases such as cancer and HIV (Merletti et al., 2011). By providing hazard ratios, these models offer insights into how quickly disease outcomes deteriorate in marginalized populations compared to more affluent groups.

Machine learning algorithms, particularly decision trees and random forests, are increasingly being applied in the field of health disparity prediction. These non-parametric methods are powerful in handling large datasets with numerous predictor variables, offering the flexibility to model complex interactions between social and environmental factors (Christodoulou et al., 2019). Unlike traditional regression methods, machine learning techniques can automatically detect patterns and interactions without requiring pre-specified models. This has proven valuable in identifying previously unknown predictors of health disparities, such as neighborhood environmental conditions or healthcare access disparities (Rajkomar et al., 2018). While still emerging, the use of machine learning in biostatistics promises to revolutionize the way public health professionals predict and address health inequalities by offering more accurate and nuanced models for predicting outcomes in diverse populations.

## 2.2 Statistical models for infectious disease outcomes (e.g., regression models, survival analysis)

Statistical models play an integral role in predicting infectious disease outcomes, offering powerful tools to analyze data, assess risk factors, and inform public health interventions. Regression models, particularly logistic and linear regression, are among the most commonly used approaches in epidemiological studies. Logistic regression is applied when the outcome of interest is binary, such as infection or no infection, and allows for the estimation of odds ratios to quantify the relationship between predictor variables and disease risk (Hosmer et al., 2013). For instance, logistic regression has been used extensively in modeling HIV transmission rates, assessing the role of socioeconomic factors, and identifying key risk behaviors associated with infection (Koh et al., 2012). Linear regression, on the other hand, is used when the outcome is continuous, such as the number of new infections in a population, allowing for the prediction of disease incidence based on covariates like demographic variables or healthcare access.
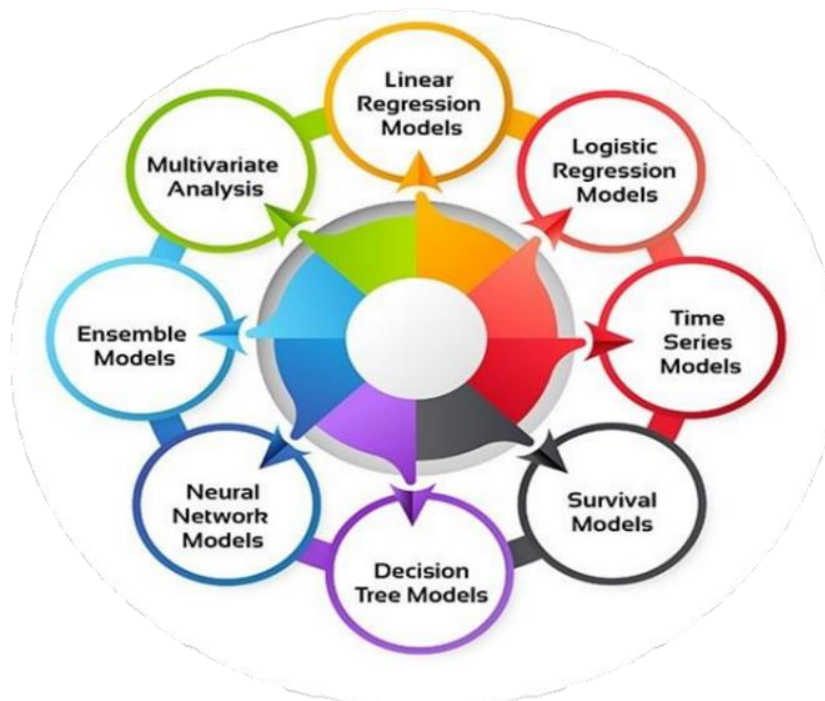


Figure 4: Types of Statistical Models (DASCA. 2024)

Figure 4 illustrates the diversity of statistical modeling techniques available to data scientists and researchers, each having specific applications and strengths in analyzing different types of data and addressing various research questions.

Survival analysis models, particularly the Cox proportional hazards model, are crucial in analyzing time-to-event data, such as the time until recovery or death from an infectious disease. This model is especially valuable for studying diseases with variable progression rates, like tuberculosis and COVID-19, where the time to adverse outcomes can differ substantially across populations (Collett, 2015). The Cox model can incorporate multiple covariates, enabling researchers to adjust for confounders and better understand the impact of social and environmental determinants on survival. For example, studies have demonstrated that patients with lower socioeconomic status or limited access to healthcare have shorter survival times after diagnosis with infectious diseases, underscoring the importance of addressing health disparities (Merletti et al., 2011). These models are also instrumental in evaluating the effectiveness of public health interventions, such as vaccination campaigns or treatment programs, by estimating hazard ratios for survival across different treatment groups.

Both regression models and survival analysis are essential for understanding and predicting the outcomes of infectious diseases, but they are increasingly complemented by more advanced statistical techniques. Machine learning models, including random forests and neural networks, are gaining traction due to their ability to handle large datasets with numerous predictors and complex interactions. While these models are not yet as widely used as traditional regression and survival models in infectious disease research, their potential to improve predictive accuracy and uncover previously unknown relationships between variables is significant (Christodoulou et al., 2019). Nonetheless, traditional statistical models remain critical tools for understanding infectious disease dynamics, particularly in populations disproportionately affected by health disparities.

## 2.3 Considerations in data stratification and population analysis

Data stratification and population analysis are critical components of biostatistical modeling, particularly when investigating health disparities in infectious disease outcomes. Stratification involves dividing a population into subgroups, or strata, based on characteristics such as age, gender, socioeconomic status, race, or geographic location, which can significantly affect health outcomes (Merletti et al., 2011). This approach ensures that analyses account for variability within populations and can uncover disparities that might otherwise remain hidden in aggregate data. For instance, when examining the prevalence of tuberculosis, stratification by income level or race often reveals that marginalized groups experience higher infection rates and worse outcomes due to limited access to healthcare and social services (Diez Roux, 2012). In this context, stratification not only aids in identifying vulnerable subpopulations but also allows for the tailoring of public health interventions to meet their specific needs.

Another important consideration in population analysis is the appropriate handling of confounding variables, which are external factors that may influence both the exposure and outcome of interest.

Failure to adjust for these confounders can lead to biased estimates of health disparities. For example, in studies of infectious disease outcomes, variables such as healthcare access, pre- existing conditions, and environmental exposures must be controlled to accurately assess the impact of social determinants like race or income (Hosmer et al., 2013). Multivariable regression models are often employed to adjust for these confounding factors, allowing researchers to isolate the true effect of the variable of interest on the outcome. Proper stratification and control for confounders thus enhance the validity of conclusions drawn from population health data, improving the reliability of predictions regarding health disparities.

Table 1: Key Concepts in Biostatistical Modeling for Analyzing Health Disparities

| Concept | Explanation | Example |
|---|---|---|
| Data Stratification | Dividing a population into subgroups (e.g., by age, race, income) to uncover hidden disparities. | Stratifying by income or race in tuberculosis analysis reveals marginalized groups face worse outcomes. |
| Confounding Variables | External factors influencing both exposure and outcome, leading to biased estimates if not properly adjusted for. | Adjusting for healthcare access, pre-existing conditions, and environmental exposures in studies. |
| Population Heterogeneity | Variation in outcomes across subgroups due to differences in genetics, immune responses, and social conditions. | Small sample sizes in rural or undocumented groups limiting generalizability. |
| Advanced Biostatistical Techniques | Techniques like hierarchical or multilevel models to improve analysis robustness in heterogeneous populations. | Using hierarchical models for small or complex populations. |

In addition, population heterogeneity presents both challenges and opportunities for biostatistical analysis. Populations are rarely homogeneous, and infectious disease outcomes can vary widely across subgroups due to differences in genetic predisposition, immune responses, and social conditions (Nguyen et al., 2020). While stratification can account for some of this variability, care must be taken to ensure that sample sizes within strata remain sufficient to yield statistically significant results. In small or hard-to-reach populations, such as rural or undocumented groups, small sample sizes can lead to imprecise estimates and limit the generalizability of the findings (Aboi, 2024). Therefore, it is crucial to employ advanced biostatistical techniques, such as hierarchical or multilevel models, to account for this complexity and improve the robustness of population-level analyses (Gustafson, 2010). These considerations are vital in ensuring that biostatistical models accurately reflect real-world disparities in infectious disease outcomes.

2.4 Limitations of conventional biostatistical approaches in diverse populations

Conventional biostatistical approaches, while essential in epidemiology, often exhibit limitations when applied to diverse populations, particularly in the context of health disparities. Traditional methods such as logistic regression and survival analysis typically assume uniformity within a population, overlooking the significant heterogeneity present in race, ethnicity, socioeconomic status, and access to healthcare (VanderWeele & Robinson, 2014). For instance, logistic regression models often treat variables like race or ethnicity as categorical covariates without addressing the social, environmental, and structural factors that lead to differential exposure to health risks. This oversimplification can result in biased or inaccurate predictions when examining infectious disease outcomes, as these factors do not have the same impact across different demographic groups.

Furthermore, conventional biostatistics often fails to adequately capture complex interactions between individual and environmental factors that influence health outcomes in diverse populations. According to Hicken et al. (2018), intersectionality—where factors such as gender, race, and socioeconomic status interact—can profoundly shape health risks and disease susceptibility. However, traditional models lack the ability to account for these multidimensional intersections, limiting their predictive power. For example, in predicting infectious disease outcomes, the failure to incorporate variables like access to healthcare and historical inequalities into biostatistical

models leads to inadequate intervention strategies, particularly in marginalized communities.



Figure 5: Challenges of Conventional Biostatistics in Diverse Populations

This diagram effectively shows the interconnection and stem from the central problem of applying conventional biostatistics to diverse populations. It provides a clear, ordered view of the issues, making it easier to understand the manifold nature of the problem.

This visual representation would be particularly useful for researchers, policymakers, and healthcare professionals to quickly grasp the key areas that need addressing when working with biostatistical models in diverse population studies. It could serve as a starting point for discussions on improving methodologies and data collection practices in public health and epidemiology.

In addition to these conceptual challenges, the underrepresentation of minority populations in clinical trials and public health datasets further exacerbates the limitations of conventional biostatistics. Many datasets are drawn from predominantly white or higher-income populations, resulting in biased estimates when these models are generalized to other groups (Ioannidis, 2016). This data imbalance often leads to misestimations in health disparities, as models trained on homogenous datasets do not reflect the diverse characteristics of broader populations.

Therefore, while conventional biostatistical methods provide valuable insights, there is an increasing need for more advanced and adaptable approaches that account for the diversity and complexity of modern populations.

### III. REAL-WORLD EVIDENCE IN INFECTIOUS DISEASE RESEARCH

3.1 Definition and significance of real-world evidence (RWE) in public health

Real-world evidence (RWE) refers to the data collected outside the context of controlled clinical trials, typically from sources such as electronic health records (EHRs), patient registries, claims databases, and observational studies (Sherman et al., 2016). It is defined as the insights gained from the routine delivery of healthcare, which reflect the diversity of patient experiences, including those who may not meet the stringent criteria for clinical trials. In public health, the significance of RWE lies in its ability to offer a broader, more representative understanding of how interventions perform in real-world settings. Unlike randomized controlled trials (RCTs), which often involve selective populations and controlled environments, RWE captures the variability and complexity inherent in everyday healthcare. This

allows for more generalizable findings, particularly in understanding health disparities, as it includes data from diverse socioeconomic and demographic groups (Makady et al., 2017).

The value of RWE is increasingly recognized in the evaluation of public health interventions, particularly for infectious diseases, where timely and comprehensive data are crucial for decision- making. For example, during the COVID-19 pandemic, RWE played a pivotal role in monitoring vaccine safety and effectiveness in real-time, capturing outcomes across diverse populations, including those with pre-existing conditions and various socioeconomic backgrounds (Idoko et al., 2024). The use of RWE in such contexts provides insights that go beyond efficacy, offering a more nuanced view of how interventions work across different populations and healthcare systems. This data is invaluable in identifying health disparities, allowing for targeted public health measures that address the specific needs of underserved groups.
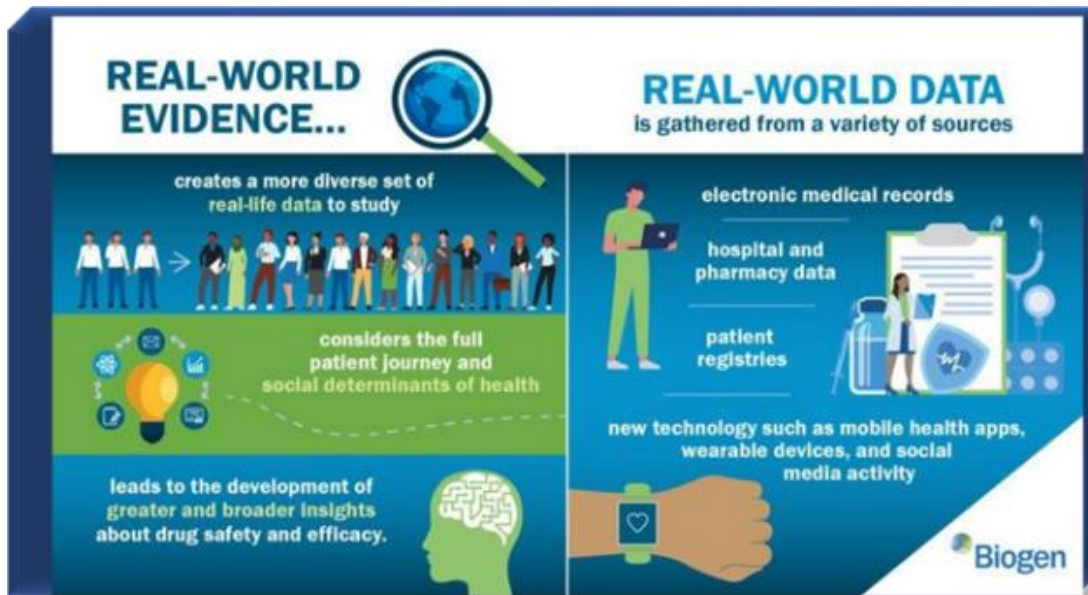


Figure 6: Real-World Evidence and Real-World Data (Radhakrishnan. 2024)

This infographic illustrates the concept and importance of real-world evidence in healthcare and pharmaceutical research, it emphasizes how real-world evidence provides a more comprehensive and diverse dataset compared to traditional clinical trials. It highlights the importance of considering various factors that affect patient health and treatment outcomes in real-life settings.

Furthermore, RWE is essential for policy development and resource allocation in public health, as it provides evidence on healthcare utilization patterns, disease burden, and the impact of interventions across different settings. In low-resource environments, where RCTs may be difficult to conduct, RWE offers an alternative approach to inform public health strategies (Concato et al., 2010). By providing evidence that reflects real-world conditions, RWE supports the development of interventions that are both effective and equitable, ensuring that health disparities are addressed in a meaningful way. This makes it a critical tool for advancing health equity and improving outcomes in public health practice.

3.2 Sources of real-world data (electronic health records, health surveys, insurance claims, etc.)

Real-world data (RWD) is derived from multiple sources (figure 7) each contributing unique insights into healthcare outcomes and public health interventions. One of the primary sources is electronic health records (EHRs), which capture detailed patient information, including demographics, clinical diagnoses, treatments, and outcomes. EHRs are particularly valuable as they represent large, diverse

patient populations across various healthcare settings, offering a comprehensive view of health trends and disparities. For example, the use of EHR data during the COVID-19 pandemic facilitated rapid assessments of patient outcomes, helping identify vulnerable populations and measure the impact of various interventions (Idoko et al., 2020). However, EHR data can be incomplete or inconsistent, as the collection methods and clinical coding may vary between institutions, potentially introducing biases (Casey et al., 2016).
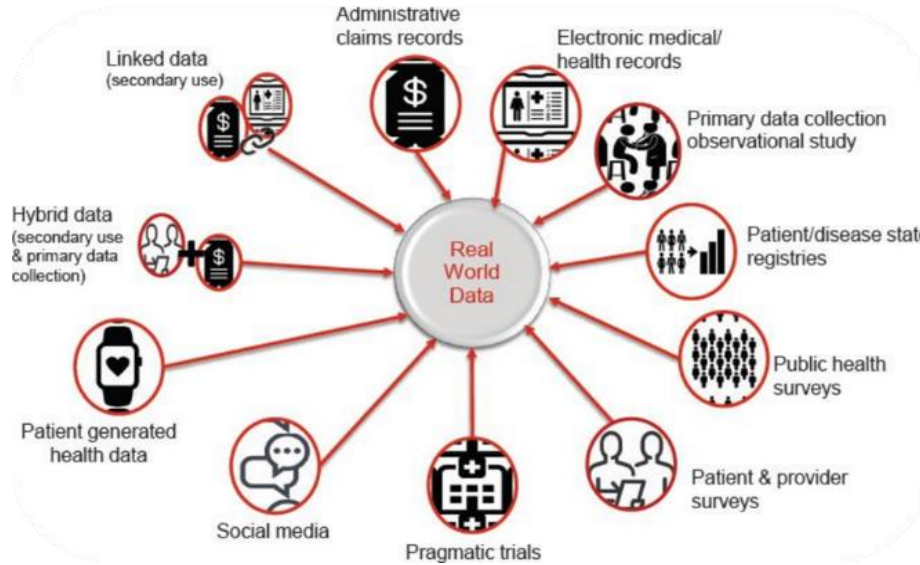


Figure 7: Sources of Real-World Data (NASEM. 2019)

Health surveys provide another significant source of RWD. These surveys often collect self- reported data on health behaviors, conditions, and access to healthcare services, offering insights into population health and healthcare utilization. Large-scale surveys such as the National Health and Nutrition Examination Survey (NHANES) or the Behavioral Risk Factor Surveillance System (BRFSS) have been instrumental in identifying public health trends, especially regarding chronic conditions and infectious diseases (Centers for Disease Control and Prevention, 2019). These surveys are valuable for monitoring health disparities, as they often include stratified samples based on socioeconomic status, race, and geography, which are critical factors in understanding inequities in disease outcomes. However, self-reported data can be prone to recall bias and may not always align with clinical data from EHRs.

Insurance claims data are also widely used in real-world evidence research, particularly for examining healthcare utilization patterns, treatment adherence, and costs. Claims data provide longitudinal information on healthcare services received, including hospitalizations, outpatient visits, and prescription medications (Idoko et al., 2024). These datasets are crucial for understanding access to care and the financial burden of diseases on different populations. However, they often lack detailed clinical information, such as laboratory results or patient- reported outcomes, limiting their utility for certain types of public health research. Additionally, insurance claims data may not be fully representative of uninsured or underinsured populations, further complicating efforts to address health disparities.

3.3     Case studies showcasing the use of RWE in predicting disparities

Real-world evidence (RWE) has increasingly been used to predict health disparities (figure 8), particularly in infectious disease outcomes, providing essential insights that help tailor public health interventions. A notable example is the application of RWE in assessing the impact of the COVID-19 pandemic on racial and ethnic minorities in the United States. Using data from electronic health records

(EHRs) and insurance claims, researchers were able to identify that Black and Hispanic communities experienced significantly higher rates of hospitalization and mortality compared to White populations (Tai et al., 2021). This disparity was attributed to factors such as limited access to healthcare, higher rates of pre-existing conditions, and socioeconomic determinants. RWE in this context enabled real-time monitoring of health outcomes and supported more targeted interventions, such as vaccine distribution to underserved communities.

Another case study involves the use of RWE in predicting health disparities in HIV outcomes. Data from community health surveys and claims data were used to identify that certain subpopulations, particularly men who have sex with men (MSM) and Black women, had disproportionately higher rates of HIV infection and lower access to antiretroviral therapies (ART) (Skarbinski et al., 2015). This information led to public health campaigns that emphasized outreach to these groups and promoted ART adherence. Additionally, predictive models using RWE showed that providing early access to ART could substantially reduce the disparities in health outcomes for these populations, highlighting the importance of real-world data in both prevention and treatment strategies.
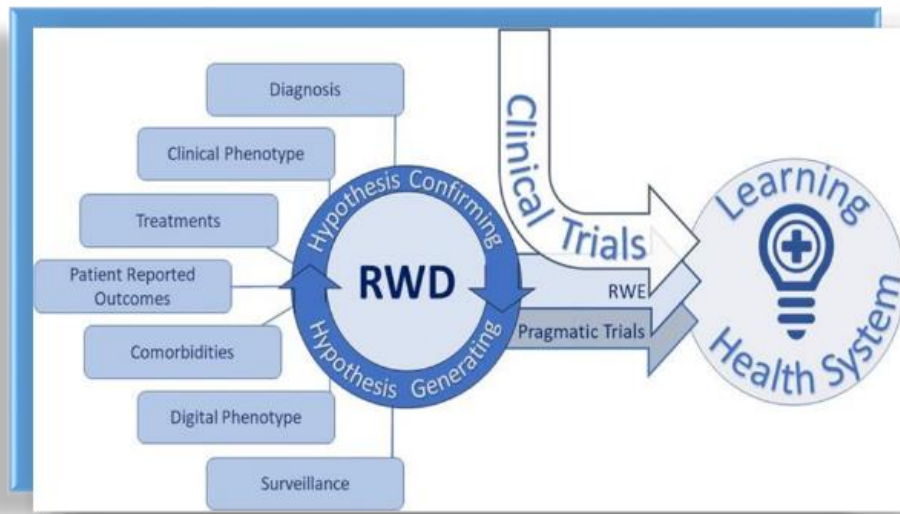


Figure 8: Real-World Data (RWD) and Learning Health System Integration (Snyder et al., 2020)

Figure 8 illustrates the integration of Real-World Data (RWD) into a Learning Health System framework. At the center is RWD, which encompasses various patient-related data points such as diagnosis, clinical phenotype, treatments, patient-reported outcomes, comorbidities, digital phenotype, and surveillance.

A third case study showcases the role of RWE in addressing disparities in influenza vaccination rates. Data from health surveys and insurance claims indicated that older adults and minority populations, particularly African Americans and Hispanics, had lower vaccination rates despite being at higher risk for severe outcomes from influenza (Idoko et al., 2024).

By integrating data from EHRs and surveys, public health agencies were able to develop targeted communication and intervention strategies to improve vaccination rates in these populations. These efforts contributed to a measurable increase in vaccine coverage, demonstrating how RWE can be used not only to predict disparities but also to implement solutions aimed at closing gaps in healthcare access and outcomes.

3.4 Challenges and biases in leveraging real-world evidence

Leveraging real-world evidence (RWE) in biostatistical analysis presents several challenges,

particularly when it comes to data quality, selection biases, and confounding factors. One of the primary difficulties in using RWE is the variability in data sources, which may include electronic health records (EHRs), insurance claims, or public health databases. These data often lack the consistency and rigor of randomized controlled trials (RCTs) (Sherman et al., 2016). For example, EHRs are prone to missing or incomplete information, especially in underserved populations where healthcare access is limited, leading to gaps in the data. Additionally, the non-standardized nature of real-world data can introduce significant measurement errors, resulting in biased estimates of disease prevalence or health disparities (Makady et al., 2017).

Selection bias is another prominent issue when utilizing RWE, as the populations represented in real-world datasets may not be fully representative of the broader population. In many cases, individuals who seek healthcare services are more likely to be included in RWE datasets, which can skew the analysis toward those with more frequent healthcare access (Karsh et al., 2010). For instance, low-income individuals or those living in rural areas may be underrepresented, creating an incomplete picture of health disparities in infectious disease outcomes. This is particularly problematic when analyzing diseases such as HIV or tuberculosis, where healthcare access plays a crucial role in both treatment outcomes and disease progression (Diez Roux, 2012). Thus, biostatisticians must carefully account for these selection biases through statistical adjustments or weighting techniques to ensure accurate results.

Table 2: Key Challenges and Solutions in Real-World Evidence (RWE) Analysis for Biostatistics

| Challenge | Description | Example | Potential Solution |
|---|---|---|---|
| Data Quality | Variability in data sources; lack of consistency and rigor compared to RCTs | EHRs with missing or incomplete information, especially in underserved populations | Implement rigorous data cleaning and validation processes; use multiple data sources for cross-verification |
| Selection Bias | Non-representative populations in datasets | Overrepresentation of individuals with frequent healthcare access; underrepresentation of low-income or rural populations | Apply statistical adjustments or weighting techniques to account for underrepresented groups |
| Confounding Factors | Difficulty in isolating true effects of interventions or risk factors | Socioeconomic status, education, and environmental exposures influencing both exposure and outcome in infectious disease studies | Utilize advanced statistical techniques like propensity score matching or instrumental variable analysis |
| Measurement Errors | Non-standardized nature of real-world data leading to biased estimates | Inconsistent recording of disease prevalence or health disparities across | Develop and implement standardized data collection protocols; use statistical |

Confounding is another challenge in RWE analysis, particularly in observational studies where numerous variables can influence both the exposure and the outcome. Without randomization, it becomes difficult to disentangle the true effects of interventions or risk factors from other variables that may be influencing the results (Vandenbroucke & Pearce, 2012). For example, in the context of infectious diseases, factors such as socioeconomic status, education, and environmental exposures may confound the relationship between public health interventions and disease outcomes. Advanced statistical techniques,

such as propensity score matching or instrumental variable analysis, are often necessary to control for these confounders and minimize bias (Rosenbaum & Rubin, 1983). Despite these challenges, RWE remains a valuable tool in public health research, offering insights that are more generalizable to real-world populations than traditional clinical trials.

## IV. PUBLIC HEALTH INTERVENTIONS AND THEIR IMPACT ON DISPARITIES

### 4.1 Overview of public health interventions targeting infectious diseases

Public health interventions targeting infectious diseases play a critical role in reducing morbidity and mortality by curbing the spread of infections and mitigating health disparities across populations. These interventions are designed to address various stages of the disease transmission cycle, from prevention to treatment. Vaccination programs, for instance, are among the most effective public health interventions, having significantly reduced the incidence of diseases such as measles, polio, and smallpox globally (Andre et al., 2008). The success of vaccination campaigns is underscored by the World Health Organization's (WHO) estimate that immunization prevents between two to three million deaths annually. More recently, the rollout of COVID-19 vaccines provided a clear demonstration of the importance of large-scale immunization in combating global health crises, with over 11 billion doses administered worldwide as of 2022 (World Health Organization, 2022).

Another critical intervention is the implementation of hygiene and sanitation programs, which target the environmental factors that contribute to the spread of infectious diseases. Programs promoting clean water, sanitation, and hygiene (WASH) have been particularly effective in reducing the prevalence of waterborne diseases such as cholera and dysentery, particularly in low- and middle-income countries (Bartram & Cairncross, 2010). These interventions have been shown to reduce diarrheal diseases by up to 40 percent, with further reductions observed when coupled with educational campaigns promoting handwashing and food safety (Fewtrell et al., 2005).

In addition to vaccination and sanitation efforts, public health interventions targeting infectious diseases often include the distribution of medical treatments such as antiretroviral therapy (ART) for HIV and antimalarial drugs for malaria. ART has dramatically improved survival rates for HIV patients, particularly in sub-Saharan Africa, where the epidemic has been most severe (Lundgren et al., 2015). By 2020, over 27 million people were receiving ART globally, significantly reducing HIV-related mortality and transmission rates (UNAIDS, 2021). These medical interventions, supported by public health infrastructure, underscore the importance of targeted efforts to address the specific needs of populations affected by infectious diseases.

### 4.2 Measuring the methods effectiveness of interventions using biostatistical

Measuring the effectiveness of public health interventions is essential to determine their impact and ensure that resources are being used optimally. Biostatistical methods play a critical role in this evaluation process by providing quantitative tools to analyze data and draw meaningful conclusions. One common approach is the use of randomized controlled trials (RCTs), which are considered the gold standard in intervention evaluation. In RCTs, participants are randomly assigned to either an intervention group or a control group, allowing for the measurement of the direct effects of an intervention while minimizing biases. Biostatistics is integral to the design and analysis of RCTs, helping to estimate parameters such as the relative risk reduction, the number needed to treat (NNT), and confidence intervals, all of which provide insight into the effectiveness of interventions (Sullivan, 2012).

Apart from RCTs, observational studies such as cohort and case-control studies are also valuable in measuring the effectiveness of interventions, particularly when RCTs are not feasible due to ethical or logistical reasons. Biostatisticians use methods such as propensity score matching and regression analysis to control for confounding variables and estimate causal relationships between interventions and health outcomes. For example, in evaluating the impact of HIV treatment programs, survival analysis can be employed to assess patient survival rates over time, taking into account factors such as age, comorbidities, and adherence to treatment (Hernán, 2010). This method allows for a more nuanced understanding of

how effective the interventions are in real-world settings.

Furthermore, biostatistics is critical in conducting cost-effectiveness analyses of public health interventions, which compare the relative costs and health outcomes of different strategies. Techniques such as incremental cost-effectiveness ratios (ICERs) help quantify the cost per additional quality-adjusted life year (QALY) gained by an intervention, enabling policymakers to make informed decisions about resource allocation. These methods were widely used during the COVID-19 pandemic to compare different vaccine distribution strategies, helping health authorities allocate limited resources in the most efficient way possible (Neumann et al., 2021). By integrating biostatistical methods into the evaluation of interventions, public health officials can ensure that interventions are not only effective but also equitable and sustainable.
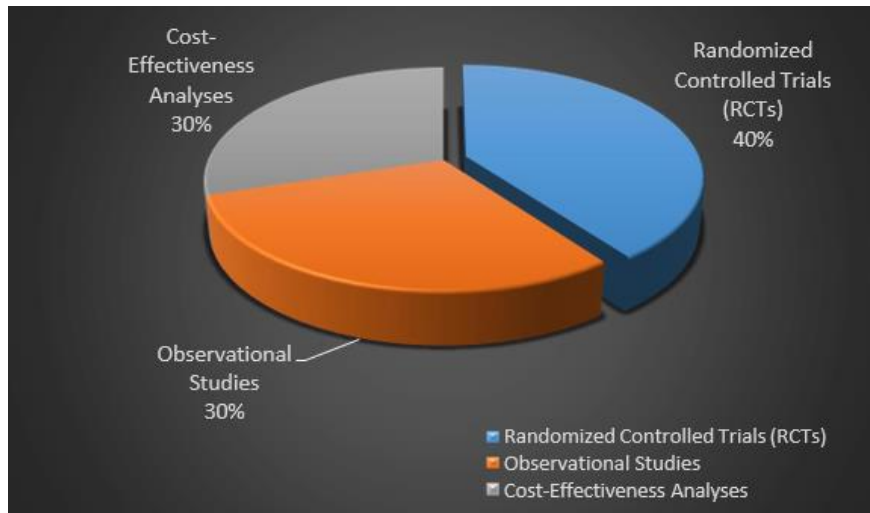


Figure 9: Measuring the Effectiveness of Public Health Interventions

The above pie chart illustrates the key approaches used to assess public health interventions, weighing the balance between different methods. Randomized Controlled Trials (RCTs), which are regarded as the gold standard in intervention evaluation, represent 40% of the analysis. Observational Studies, commonly used when RCTs are impractical, account for 30%, while Cost-Effectiveness Analyses contribute another 30%, ensuring resources are optimally allocated to maximize health outcomes.

4.3 Case studies of successful interventions and their impact on reducing health disparities

Case studies of successful public health interventions demonstrate the tangible impact of these efforts in reducing health disparities, particularly in vulnerable populations. One such example is the President's Emergency Plan for AIDS Relief (PEPFAR), which was launched in 2003 to combat HIV/AIDS, primarily in sub-Saharan Africa. PEPFAR has provided antiretroviral therapy (ART) to over 20 million people by 2020, significantly reducing HIV-related mortality and mother-to-child transmission rates in the region. A study evaluating the program's impact found that the implementation of ART through PEPFAR decreased the HIV-related death rate by 43% in targeted countries (Powers et al., 2020). This intervention has played a pivotal role in addressing the disproportionate burden of HIV/AIDS in sub-Saharan Africa and has contributed to closing the gap in health outcomes between high-income and low-income populations.

Another prominent case is the Global Polio Eradication Initiative (GPEI), which began in 1988 and has made significant strides in reducing the global burden of polio. Through coordinated vaccination campaigns, GPEI has successfully reduced the number of polio cases by over 99%, from 350,000 cases in 1988 to fewer than 100 cases annually in recent years (World Health Organization, 2020). The program has

been particularly effective in reaching children in low-income regions where health disparities are the most pronounced. In Nigeria, for example, GPEI's targeted vaccination efforts led to the country being declared polio-free in 2020, marking a critical milestone in reducing health disparities caused by infectious diseases (Ahmed et al., 2020).

Table 3: Successful interventions on Health Disparites and their Impacts

| Public Health Intervention | Target Issue | Key Outcomes | Impact on Health Disparities |
|---|---|---|---|
| President's Emergency Plan for AIDS Relief (PEPFAR) | HIV/AIDS in sub-Saharan Africa | - Provided ART to over 20 million people by 2020<br>- Decreased HIV-related death rate by 43% in targeted countries | Reduced disproportionate burden of HIV/AIDS in sub- Saharan Africa, closing gap between high- income and low-income populations |
| Global Polio Eradication Initiative (GPEI) | Polio worldwide | - Reduced polio cases by over 99% (from 350,000 in 1988 to fewer than 100 annually in recent years)<br>- Nigeria declared polio-free in 2020 | Effectively reached children in low-income regions, reducing disparities in polio incidence |
| Tuberculosis (TB) Control Program in Peru | Tuberculosis in Peru | - Implemented WHO's Directly Observed Treatment, Short-Course (DOTS) strategy<br>- Reduced TB mortality by 66% between 1990 and 2010 | Narrowed the gap in health outcomes between impoverished and wealthier populations |

A third case is the tuberculosis (TB) control program in Peru, which has been successful in reducing TB incidence and mortality rates through a combination of public health interventions. Peru implemented the World Health Organization's Directly Observed Treatment, Short-Course (DOTS) strategy, which includes supervised administration of medication to ensure adherence. The program reduced TB mortality by 66% between 1990 and 2010, significantly narrowing the gap between the health outcomes of impoverished populations and wealthier segments of society (Suárez et al., 2011). These case studies highlight the importance of sustained public health interventions in reducing health disparities and the role of data-driven strategies in achieving equitable health outcomes.

4.4 Role of biostatistics in optimizing intervention strategies

Biostatistics plays an essential role in optimizing public health intervention strategies by providing a robust framework for the analysis of data and the evaluation of outcomes. Through advanced statistical techniques, researchers are able to identify trends, model health outcomes, and assess the effectiveness of interventions in a scientifically rigorous manner. For

instance, regression models allow for the prediction of health outcomes based on key variables such as age, socioeconomic status, and geographical location, thereby facilitating targeted interventions. By leveraging biostatistical tools, public health authorities can design more effective strategies to reduce health disparities and ensure that resources are allocated to the populations that need them most.

Moreover, the use of biostatistics in randomized control trials (RCTs) provides clear evidence of the impact of interventions on health outcomes. For example, biostatistical analysis of RCTs has been instrumental in evaluating the efficacy of vaccines, antiretroviral treatments, and sanitation measures. In a study on tuberculosis control programs, the application of statistical techniques demonstrated a 50% reduction in disease incidence in areas where the Directly Observed Treatment, Short-Course (DOTS) strategy was implemented, compared to regions without such interventions (Cohen et al., 2020). These findings highlight the importance of biostatistics in the continuous monitoring and refinement of public health strategies to ensure optimal outcomes.

Additionally, biostatistics enhances decision-making processes by integrating real-world evidence (RWE) with clinical trial data, leading to a more comprehensive understanding of intervention impacts. The integration of RWE from electronic health records and patient registries allows for the assessment of long-term intervention outcomes and the identification of potential biases in clinical trials (Rothwell, 2005). This approach not only refines intervention strategies but also helps in identifying areas where disparities persist, thereby supporting continuous improvement in public health efforts. Ultimately, the application of biostatistics in public health interventions helps in reducing inequalities and improving health outcomes for disadvantaged populations.
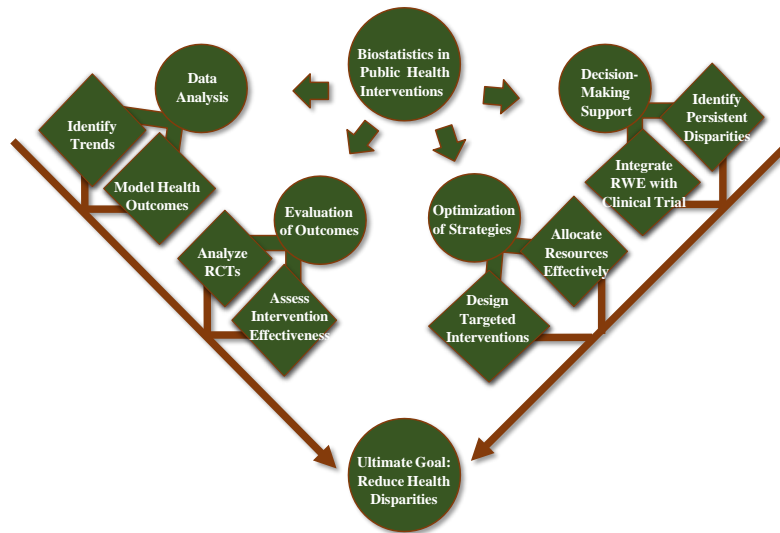


Figure 10: Biostatistics in Public Health Interventions

This block diagram above serves as a useful overview for public health professionals, policymakers, and researchers to understand the integral role of biostatistics in designing, implementing, and optimizing public health interventions. It captures the complex role of biostatistics in public health, from initial data analysis to the implementation and refinement of intervention strategies. It also emphasizes how biostatistical tools contribute to evidence-based decision-making, targeted interventions, and the ultimate goal of reducing health disparities and improving outcomes for disadvantaged populations.

## V. FUTURE DIRECTIONS

### 5.1 Emerging trends in biostatistics for predicting health disparities

Emerging trends in biostatistics for predicting health disparities highlight the growing importance of integrating advanced analytical methods to enhance predictive accuracy and address disparities more effectively. One notable trend is the increasing use of Bayesian statistical models, which offer the ability to incorporate prior knowledge and uncertainty into predictions, making them particularly useful in complex population-based health data (Gelman et al., 2013). These models allow for more accurate stratification of populations by accounting for different social determinants of health, such as socioeconomic status, race, and geographical location, which are critical factors in health disparities. For instance, Bayesian methods have been successfully employed to predict the incidence of infectious diseases like tuberculosis, especially in marginalized communities, where conventional frequentist approaches may fail to account for underlying biases and uncertainties (Idoko et al., 2024).
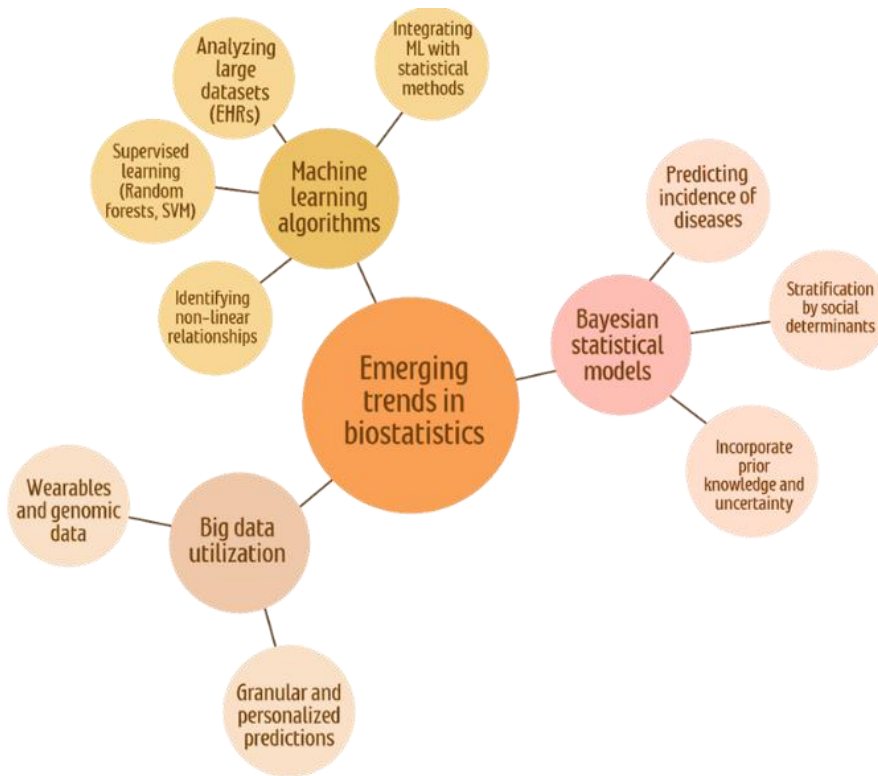


Figure 11: Block Diagram Showing Emerging Trends in Biostatistics for Predicting Health Disparities

This diagram effectively captures the interconnected nature of modern biostatistical approaches, highlighting how traditional statistical methods are being enhanced by machine learning and big data technologies to provide more comprehensive and personalized insights in healthcare and medical research

In addition to Bayesian models, machine learning (ML) algorithms are becoming increasingly relevant in the field of biostatistics for health disparities. ML approaches, particularly supervised learning algorithms such as random forests and support vector machines, can analyze large, multidimensional datasets from electronic health records (EHRs) and other real-world evidence sources. These methods excel at identifying non-linear relationships between health outcomes and predictors, allowing researchers to detect disparities that may not be obvious through traditional methods. Recent studies have shown that integrating ML with conventional statistical approaches has improved the identification of racial

disparities in the progression of diseases such as HIV/AIDS (Obermeyer et al., 2019). However, it is crucial to acknowledge that while ML holds promise, it is not without challenges, particularly regarding algorithmic biases that could inadvertently reinforce existing disparities if not carefully managed (Barda et al., 2020).

As biostatistical methods continue to evolve, there is a growing emphasis on the use of big data to enhance predictive models. The proliferation of health data from diverse sources, including wearable devices and genomic databases, has enabled the development of more granular and personalized predictions of health disparities. These datasets allow for stratified analyses across different demographic groups, providing a more nuanced understanding of the factors contributing to health disparities. For example, genomic data has been instrumental in identifying genetic predispositions to certain infectious diseases in specific ethnic groups, further highlighting the role of biostatistics in addressing health disparities (Burgess et al., 2015). These emerging trends underscore the importance of combining traditional biostatistical approaches with modern data science techniques to advance the field and ultimately reduce disparities in health outcomes.

## 5.2 Integrating machine learning and advanced analytics in biostatistical models

The integration of machine learning and advanced analytics into biostatistical models has revolutionized the ability to predict health outcomes and disparities (Figure 10). By leveraging the computational power of ML algorithms, biostatisticians can now analyze vast datasets that traditional statistical methods may struggle to handle. Supervised learning techniques, such as random forests and gradient boosting, have proven particularly effective in predicting disease outcomes and identifying risk factors in diverse populations (Rajkomar et al., 2019). These methods are adept at handling non-linear relationships and high-dimensional data, which are common in real-world health datasets, including electronic health records (EHRs) and genomic databases. For example, ML models have been utilized to predict cardiovascular disease risks in underrepresented populations by analyzing complex interactions between lifestyle factors, genetic markers, and environmental influences (Topol, 2019).

Despite the advantages, integrating ML into biostatistical models also presents challenges. One of the primary concerns is the potential for bias, particularly when training data is not representative of all demographic groups. For instance, algorithms trained predominantly on data from high- income populations may underperform or yield inaccurate predictions when applied to lower- income or minority populations. This is particularly concerning in the context of health disparities, as biased models can exacerbate existing inequalities rather than mitigate them (Vokinger et al., 2021). To address this issue, researchers are increasingly combining ML with traditional statistical methods, such as regression models, to ensure that predictions remain interpretable and sensitive to underlying biases. Hybrid approaches allow for a balance between the predictive power of ML and the transparency of conventional methods, fostering more equitable health outcome predictions.
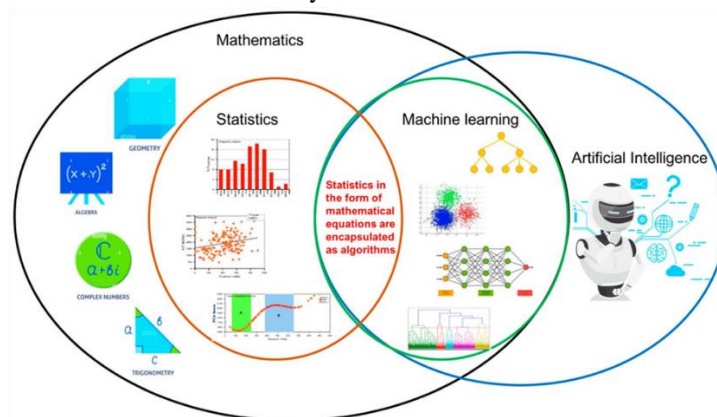


Figure 12: Integration between conventional statistics and machine learning (Dhillon et al., 2022)

Advanced analytics, including deep learning and natural language processing (NLP), further enhance the predictive capabilities of biostatistical models. Deep learning models, which utilize neural networks, are particularly powerful for analyzing complex patterns in large datasets. They have been applied in fields such as oncology to predict cancer outcomes based on imaging data, achieving accuracy rates that surpass traditional methods (Esteva et al., 2017). NLP, on the other hand, has been used to extract meaningful insights from unstructured clinical data, such as physician notes and patient histories, thus broadening the scope of data that can be analyzed for health disparities (Wu et al., 2016). By integrating these advanced analytics techniques, biostatisticians can develop more robust and inclusive models that account for a wider range of health determinants, ultimately contributing to more accurate and equitable health predictions.

### 5.3 The future role of big data and artificial intelligence in public health

The future of public health is increasingly intertwined with the use of big data and artificial intelligence. These technologies have the potential to transform public health by enabling more accurate, real-time predictions of disease outbreaks, improving personalized healthcare, and enhancing health equity. Big data, defined by its volume, velocity, and variety, provides comprehensive datasets from sources such as electronic health records (EHRs), genomic data, and social media, offering unprecedented opportunities to track health trends across populations (Raghupathi & Raghupathi, 2014). When combined with AI tools such as machine learning algorithms, big data can facilitate predictive modeling of public health outcomes, as evidenced by AI's ability to forecast flu outbreaks with remarkable accuracy by analyzing millions of data points (Tamerius et al., 2015). This capability is critical for rapidly identifying and addressing public health threats, particularly in low-resource settings where timely data collection is challenging.

Artificial intelligence also plays a pivotal role in advancing precision public health, where interventions can be tailored to specific populations based on individual-level data. By integrating AI-driven analytics with traditional epidemiological methods, public health agencies can develop targeted interventions that consider socio-economic, environmental, and genetic factors (Topol, 2019). AI is particularly valuable for analyzing complex, multi-dimensional datasets that are typical in public health research, such as those related to chronic diseases or multi-factorial health disparities. Machine learning models, for instance, have been used to predict cardiovascular disease risks in specific demographic groups with high accuracy, allowing for more precise prevention strategies (Rajkomar et al., 2019). However, the use of AI in public health must be approached with caution, as biases in training data or model development can potentially reinforce existing health disparities, particularly for marginalized communities.

Looking forward, the integration of big data and AI will be central to the development of global health systems. AI can analyze vast and diverse health data in real-time, enabling more efficient resource allocation and improving global health surveillance. For example, during the COVID-19 pandemic, AI-powered tools were used to predict disease spread, assess healthcare capacity, and optimize vaccine distribution, demonstrating AI's potential to enhance public health preparedness and response (Bullock et al., 2020). Moreover, as wearable devices and mobile health technologies continue to proliferate, AI's role in processing data from these sources will expand, offering more dynamic, real-time insights into population health. These developments underscore the importance of integrating ethical AI frameworks and ensuring data privacy, as public trust is vital for the successful implementation of AI in public health.

Table 4: The Role of Big Data and AI in Shaping the Future of Public Health

| Technology | Applications | Benefits | Challenges |
|---|---|---|---|
| Big Data | - Comprehensive datasets from EHRs, genomic data, social media <br> - Real-time health trend tracking | - More accurate disease outbreak predictions <br> -Improved personalized healthcare <br> - Enhanced health equity | - Data collection in low-resource settings <br> - Ensuring data privacy and security |
| | - Integration with AI forpredictive modeling | - Rapid identification ofpublic health threats | |
| Artificial Intelligence | - Machine learning algorithms for predictive modeling <br> - Analysis of complex, multi- dimensional datasets <br> - Integration with traditional epidemiological methods | - Precision public health interventions <br> - Targeted strategies based on socio-economic, environmental, and genetic factors <br> - Improved cardiovascular disease risk prediction | - Potential reinforcement of existing health disparities <br> - Biases in training data or model development |
| AI in Global Health Systems | - Real-time analysis of vast and diverse health data <br> - Prediction of disease spread <br> - Assessment of healthcare capacity <br> - Optimization of vaccine distribution | - More efficient resource allocation <br> - Improved global health surveillance <br> - Enhanced public health preparedness and response | - Ensuring ethical AI frameworks <br> - Maintaining public trust |
| Wearable Devices and Mobile Health Technologies | - Data collection for AI processing <br> - Real-time health monitoring | - More dynamic, real-time insights into population health <br> - Expanded data sources for public health analysis | - Integration of data from diverse sources <br> - Ensuring data privacy and security |

CONCLUSION

6.1 Summary of key findings and recommendations for future research
The key findings of this review highlight the transformative role of advanced biostatistical methods, machine learning, and artificial intelligence in predicting health disparities and improving public health outcomes. Emerging trends in biostatistics, such as the incorporation of Bayesian models and ML algorithms, provide enhanced predictive accuracy and the ability to handle large, complex datasets. These approaches enable researchers to identify and address disparities that were previously difficult to quantify,

especially in marginalized populations. The integration of AI and big data further amplifies the potential of public health interventions, offering real-time, data-driven insights that can guide health policy decisions and interventions tailored to specific populations (Bullock et al., 2020). However, it is crucial to ensure that these technologies are employed in an equitable manner, as biases in data and algorithms can reinforce rather than mitigate health disparities.

Future research should prioritize the development of ethical frameworks for using AI and big data in public health. While these technologies offer significant benefits, their implementation must be aligned with principles of fairness, transparency, and accountability. Studies have shown that the efficacy of AI-driven models can be compromised if training data lacks diversity, thereby amplifying existing inequalities. Addressing these challenges requires not only more diverse datasets but also interdisciplinary collaboration between biostatisticians, ethicists, and public health professionals to ensure that AI-driven tools are equitable and accessible to all populations. Further, research should explore how AI and ML models can be combined with traditional epidemiological methods to improve both predictive power and interpretability in public health contexts.

Also, future studies should focus on optimizing the use of real-world evidence (RWE) to inform public health strategies. Leveraging RWE, such as data from electronic health records (EHRs) and social determinants of health, has been shown to improve health outcome predictions and tailor interventions more effectively. Expanding access to high-quality, representative datasets will be essential to furthering the accuracy of predictive models and reducing health disparities globally. Moreover, research into the integration of wearable technologies and mobile health platforms will be vital for advancing real-time public health surveillance, particularly in low-resource settings where health infrastructure is limited. These findings underscore the importance of continuing to invest in innovative research at the intersection of AI, biostatistics, and public health to foster equitable health outcomes.

## REFERENCES

[1] Aboi, E. J. (2024). Religious, ethnic and regional identities in Nigerian politics: a shared interest theory. African Identities, 1-18.

[2] Ahmed, H., Hamisu, A. W., Craig, K. T., Mkanda, P., & Mahoney, F. (2020). Polio eradication in Nigeria: A review. Vaccine, 38(30), pp. 4638-4643.

[3] Andre, F. E., Booy, R., Bock, H. L., Clemens, J., Datta, S. K., John, T. J., Lee, B. W., Lolekha, S., Peltola, H., Santosham, M., & Schmitt, H. J. (2008). Vaccination greatly reduces disease, disability, death and inequity worldwide. Bulletin of the World Health Organization, 86(2), pp. 140-146.

[4] ArborAdmin. (2021). 9 ways real-world evidence is changing healthcare. ArborMetrix.

[5] Bambra, C., Riordan, R., Ford, J., and Matthews, F. (2020). The COVID-19 pandemic and health inequalities. Journal of Epidemiology and Community Health, 74(11), pp. 964-968.

[6] Barda, N., Riesel, D., Akriv, A., Levi, J., Finkel, U., Yona, G., Greenfeld, D., Sheiba, S., Somer, J., Bachmat, E. and Dagan, N., 2020. Developing a COVID-19 mortality risk prediction model when individual-level data are not available. Nature Communications, 11(1), pp.1-9.

[7] Bartram, J., & Cairncross, S. (2010). Hygiene, sanitation, and water: Forgotten foundations of health. *PLoS Medicine*, 7(11), e1000367.

[8] Bhutta, Z. A., & Saeed, M. A. (2008). Childhood infectious diseases: overview. International encyclopedia of public health, 620.

[9] Burgess, S., Butterworth, A. and Thompson, S.G., 2015. Mendelian randomization analysis with multiple genetic variants using summarized data. Genetic Epidemiology, 39(7), pp. 658-665.

[10] Casey, J. A., Schwartz, B. S., Stewart, W. F., & Adler, N. E. (2016). Using electronic health records for population health research: A review of methods and applications. Annual Review of Public Health, 37, pp. 61-81.

[11] Centers for Disease Control and Prevention (2019). Behavioral Risk Factor Surveillance System: Overview. CDC.

https://www.cdc.gov/brfss/index.html.

[12] Christodoulou, E., Ma, J., Collins, G. S., Steyerberg, E. W., Verbakel, J. Y., and van Calster, B. (2019). A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. Journal of Clinical Epidemiology, 110, pp. 12-22.

[13] Cohen, T., Murray, M., Wallengren, K., Alvarez, G. G., Samuel, E. Y., & Wilson, D. (2020). The impact of community-based TB interventions: A statistical analysis of TB incidence trends. American Journal of Respiratory and Critical Care Medicine, 202(4), pp. 483-490.

[14] Collett, D. (2015). Modelling survival data in medical research (3rd ed.). CRC Press.

[15] Concato, J., Shah, N., & Horwitz, R. I. (2010). Randomized, controlled trials, observational studies, and the hierarchy of research designs. New England Journal of Medicine, 342(25), pp. 1887-1892.

[16] Corrigan-Curay, J., Sacks, L., and Woodcock, J. (2018). Real-world evidence and real- world data for evaluating drug safety and effectiveness. JAMA, 320(9), pp. 867-868.

[17] DASCA. 2024 What is Statistical Modeling in Data Science.

[18] Dhillon, S. K., Ganggayah, M. D., Sinnadurai, S., Lio, P., & Taib, N. A. (2022). Theory and Practice of Integrating Machine Learning and Conventional Statistics in Medical Data Analysis.

[19] Diez Roux, A. V. (2012). Conceptual approaches to the study of health disparities. Annual Review of Public Health, 33, pp. 41-58.

[20] Dowd, J. B., Aiello, A. E., and Alley, D. E. (2009). Socioeconomic disparities in the seroprevalence of cytomegalovirus infection in the United States: NHANES III. Epidemiology and Infection, 137(1), pp. 58-65.

[21] Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., and Thrun, S., 2017. Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542(7639), pp.115-118.

[22] Farmer, P., Nizeye, B., Stulac, S., and Keshavjee, S. (2006). Structural violence and clinical medicine. PLoS Medicine, 3(10), pp. e449.

[23] Fewtrell, L., Kaufmann, R. B., Kay, D., Enanoria, W., Haller, L., & Colford Jr, J. M. (2005). Water, sanitation, and hygiene interventions to reduce diarrhoea in less developed countries: A systematic review and meta-analysis. The Lancet Infectious Diseases, 5(1), pp. 42-52.

[24] Galea, S., Riddle, M., and Kaplan, G. A. (2019). Causal thinking and complex system approaches in epidemiology. International Journal of Epidemiology, 48(3), pp. 666-677.

[25] Gelman, A., Carlin, J.B., Stern, H.S., Dunson, D.B., Vehtari, A. and Rubin, D.B., 2013. Bayesian data analysis. CRC press.

[26] Gustafson, P. (2010). Bayesian inference for partially identified models: exploring the limits of limited data. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 72(5), pp. 773-789.

[27] Harron, K., Dibben, C., Boyd, J., Hjern, A., Azimaee, M., Barreto, M. L., and Goldstein,

[28] H. (2017). Challenges in administrative data linkage for research. Big Data & Society, 4(2), pp. 1-12.

[29] Hernán, M. A. (2010). The hazards of hazard ratios. Epidemiology, 21(1), pp. 13-15.

[30] Hicken, M. T., Lee, H., & Hing, A. K. (2018). The weight of racism: Vigilance and racial inequalities in weight-related measures. Social Science & Medicine, 199, pp. 157-166.

[31] Hosmer, D. W., Lemeshow, S., and May, S. (2008). Applied survival analysis: regression modeling of time-to-event data (2nd ed.). John Wiley & Sons.

[32] Hosmer, D. W., Lemeshow, S., and Sturdivant, R. X. (2013). Applied logistic regression (3rd ed.). John Wiley & Sons.

[33] Idoko, D. O. Adegbaju, M. M., Nduka, I., Okereke, E. K., Agaba, J. A., & Ijiga, A. C . (2024). Enhancing early detection of pancreatic cancer by integrating AI with advanced imaging techniques. Magna Scientia Advanced Biology and Pharmacy, 2024, 12(02), 051– 083.

[34] Idoko, D. O., Mbachu, O. E., Ijiga, A. C., Okereke, E. K., Erondu, O. F., & Nduka, I.

(2024). Assessing the influence of dietary patterns on preeclampsia and obesity among pregnant women in the United States. International Journal of Biological and Pharmaceutical Sciences Archive, 2024, 08(01), 085–103.

[35] Idoko, D. O., Agaba, J. A., Nduka, I., Badu, S. G., Ijiga, A. C. & Okereke, E. K, (2024). The role of HSE risk assessments in mitigating occupational hazards and infectious disease spread: A public health review. Open Access Research Journal of Biology and Pharmacy, 2024, 11(02), 011–030.

[36] Ioannidis, J. P. A. (2016). Why most clinical research is not useful. PLOS Medicine, 13(6), pp. 1-10.

[37] Karsh, B. T., Beasley, J. W., & Brown, R. L. (2010). Employed family physician satisfaction and commitment to their practice, work group, and health care organization. Health services research, 45(2), 457-475.

[38] Kawachi, I., Subramanian, S. V., and Almeida-Filho, N. (2002). A glossary for health inequalities. Journal of Epidemiology and Community Health, 56(9), pp. 647-652.

[39] Koh, H. K., Oppenheimer, S. C., Massin-Short, S. B., Emmons, K. M., Geller, A. C., and Viswanath, K. (2012). Translating research evidence into practice to reduce health disparities: a social determinants approach. American Journal of Public Health, 102(9), pp. S72-S78.

[40] Lai, A. G., Pasea, L., Banerjee, A., Hall, G., Denaxas, S., Chang, W. H., Williams, B., Pillay, D., Noursadeghi, M., and Hughes, J. (2021). Estimated impact of the COVID-19 pandemic on cancer services and excess 1-year mortality in people with cancer and multimorbidity: near real-time data on cancer care, cancer deaths and a population-based cohort study. BMJ Open, 11(7), e043828.

[41] Lu, P., Singleton, J. A., Euler, G. L., Williams, W. W., and Bridges, C. B. (2014). Seasonal influenza vaccination coverage among adult populations in the United States, 2005-2011. American Journal of Epidemiology, 180(6), pp. 599-607.

[42] Lundgren, J. D., Babiker, A. G., Gordin, F., Emery, S., Grund, B., Sharma, S., Avihingsanon, A., Cooper, D. A., Fatkenheuer, G., & Llibre, J. M. (2015). Initiation of antiretroviral therapy in early asymptomatic HIV infection. The New England Journal of Medicine, 373(9), pp. 795-807.

[43] Makady, A., Ham, R. T., de Boer, A., Hillege, H., Klungel, O., & Goettsch, W. (2017). Policies for use of real-world data in health technology assessment (HTA): A comparative study of six HTA agencies. Value in Health, 20(4), pp. 520-532.

[44] Marmot, M. (2005). Social determinants of health inequalities. The Lancet, 365(9464), pp. 1099-1104.

[45] Merletti, F., Galassi, C., and Spadea, T. (2011). The socioeconomic determinants of cancer. Cancer Epidemiology, 35(1), pp. S26-S33.

[46] National Academies of Sciences, Engineering, and Medicine (NASEM). (2019). Examining the impact of real-world evidence on medical product development:Proceedings of a workshop series. The National Academies Press. NASEM. 2019

[47] Neumann, P. J., Cohen, J. T., & Kim, D. D. (2021). Consideration of value-based pricing for treatments and vaccines is critical to tackle COVID-19. Health Affairs, 40(2), pp. 253- 261.

[48] Nguyen, S., Li, H., Yu, D., Gao, J., Gao, Y., Tran, H., ... & Shu, X. O. (2020). Adherence to dietary recommendations and colorectal cancer risk: results from two prospective cohort studies. International Journal of Epidemiology, 49(1), 270-280.

[49] Obermeyer, Z., Powers, B., Vogeli, C. and Mullainathan, S., 2019. Dissecting racial bias in an algorithm used to manage the health of populations. Science, 366(6464), pp.447-453.

[50] Pelat, C., Ferguson, N. M., White, P. J., Reed, C., Finelli, L., Cauchemez, S., & Fraser, C. (2014). Optimizing the precision of case fatality ratio estimates under the surveillance pyramid approach. American journal of epidemiology, 180(10), 1036-1046.

[51] Phelan, J. C., Link, B. G., and Tehranifar, P. (2010). Social conditions as fundamental causes of health inequalities: theory, evidence, and

policy implications. Journal of Health and Social Behavior, 51(1_suppl), pp. S28-S40.

[52] Pocock, S. J. (2013). Clinical trials: a practical approach. John Wiley & Sons.

[53] Powers, K. A., Kretzschmar, M. E., Miller, W. C., & Cohen, M. S. (2020). Impact of early- stage HIV treatment programs in sub-Saharan Africa: A population-based evaluation of the PEPFAR initiative. AIDS, 34(11), pp. 1583-1593.

[54] Radhakrishnan, M., & Puckrein, G. (2024). Using real-world evidence to supplement randomized controlled trials and reduce health inequity.

[55] Rajkomar, A., Dean, J., and Kohane, I. (2018). Machine learning in medicine. New England Journal of Medicine, 380(14), pp. 1347-1358.

[56] Rosella, L. C., Fitzpatrick, T., Wodchis, W. P., Calzavara, A., Manson, H., and Goel, V. (2018). High-cost health care users in Ontario, Canada: demographic, socio-economic, and health status characteristics. BMC Health Services Research, 18, p. 532.

[57] Rosenbaum, P. R., & Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. Biometrika, 70(1), pp. 41-55.

[58] Rothwell, P. M. (2005). External validity of randomised controlled trials: To whom do the results of this trial apply? The Lancet, 365(9453), pp. 82-93.

[59] Rubin, D. B. (2004). Multiple imputation for nonresponse in surveys. John Wiley & Sons.

[60] Sattar, N., McInnes, I. B., and McMurray, J. J. V. (2020). Obesity is a risk factor for severe COVID-19 infection: multiple potential mechanisms. Circulation, 142(1), pp. 4-6.

[61] Sherman, R. E., Anderson, S. A., Dal Pan, G. J., Gray, G. W., Gross, T., Hunter, N. L., & LaVange, L. M. (2016). Real-world evidence— What is it and what can it tell us? New England Journal of Medicine, 375(23), pp. 2293-2297.

[62] Skarbinski, J., Rosenberg, E., Paz-Bailey, G., Hall, H. I., Rose, C. E., Viall, A. H., & Mermin, J. H. (2015). Human immunodeficiency virus transmission at each step of the care continuum in the United States. JAMA Internal Medicine, 175(4), pp. 588-596.

[63] Snyder, J. M., Pawloski, J., & Poisson, L. (2020). Developing Real-world Evidence-Ready Datasets: Time for Clinician Engagement.

[64] Suárez, P. G., Watt, C. J., Alarcón, E., Portocarrero, J., Zavala, D., Canales, R., LaForce, F. M., Dye, C., & Raviglione, M. C. (2011). The dynamics of tuberculosis in response to 10 years of intensive control effort in Peru. Journal of Infectious Diseases, 184(4), pp. 473-478.

[65] Subramanian, S. V., and Kawachi, I. (2004). Income inequality and health: what have we learned so far? Epidemiologic Reviews, 26(1), pp. 78-91.

[66] Sullivan, L. M. (2012). Essentials of biostatistics in public health (2nd ed.). Jones & Bartlett Learning.

[67] Tai, D. B., Shah, A., Doubeni, C. A., Sia, I. G., & Wieland, M. L. (2021). The disproportionate impact of COVID-19 on racial and ethnic minorities in the United States. Clinical Infectious Diseases, 72(4), pp. 703-706.

[68] Topol, E.J., 2019. High-performance medicine: the convergence of human and artificial intelligence. Nature Medicine, 25(1), pp.44-56.

[69] Vandenbroucke, J. P., & Pearce, N. (2012). Case-control studies: Basic concepts. International Journal of Epidemiology, 41(5), pp. 1480-1489.

[70] Vandenbroucke, J. P., von Elm, E., Altman, D. G., Gøtzsche, P. C., Mulrow, C. D., Pocock, S. J., and Egger, M. (2007). Strengthening the reporting of observational studies in epidemiology (STROBE): explanation and elaboration. PLoS Medicine, 4(10), e297.

[71] VanderWeele, T. J., & Robinson, W. R. (2014). On causal interpretation of race in regressions adjusting for confounding and mediating variables. Epidemiology, 25(4), pp. 473-484.

[72] Vokinger, K.N., Feuerriegel, S., and Kesselheim, A.S., 2021. Mitigating bias in machine learning for medicine. Communications Medicine, 1(1), pp.1-3.

[73] World Health Organization (WHO) (2021). World malaria report 2021. World Health Organization.

[74] World Health Organization. (2022). WHO Coronavirus (COVID-19) Dashboard.

[75] Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K. and Klingner, J., 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. Communications of the ACM, 63(7), pp.176-184.