**DeepScience**
Open Access Books

# Chapter 8: The role of big data analytics in understanding population health and social determinants of disease

## 8.1. Introduction

Many organizations, systems, businesses, and institutions are inundated with ever-increasing amounts of data and information. In the health sector, monitoring changes in population health data can provide guidance on resource management, policy adjustment, and disease control. This article analyzes and discusses examples of big health data analytics, and some needs and directions for improving the use of big data for health analytics in the context of the currently ongoing epidemic situation in China. With the widespread use of the Internet, a huge amount of health data is generated. Of particular note are several main types of health data: hospital health records, electronic health records, health insurance claim records, medical records and results of clinical trials. Most health data are now being stored in electronic form, making it accessible and facilitation. Health data are more complex than many other types of data due to their privacy and relevance.

Software, hardware, and network infrastructures are integrated with public health practice and data resources to strengthen information processing and storage capabilities. Active cooperation between academia, scientific research labs, and government regulatory agencies could foster the development of databases and data standards for emerging areas. The introduction of new models of data sharing and collaboration, supported by strong privacy and safety measures, could stimulate valuable analytic research and maximize the social utility of research data. Public Health Potential Realization in Big Data Era: Expansion of the conceptual domain of public health to include other public sector and private sector actors influencing population health outcomes; The integration of traditional public health and big data analytics frameworks through systematic attention to context; The adoption by Big Data of traditional public

health values, including attention to actionable knowledge, equity, and the importance of empirical validation in the knowledge-production process.

The introduction of the Big Data concept in the healthcare sector indicates a great challenge and potential. The Big data analytics solutions are developing rapidly and are greatly practical. The goal is to indicate the importance of analyzing and processing large volumes of data that go beyond the typical ways of storing and processing information in the healthcare system.



**Fig 8.1:** Big data analytics in healthcare

### 8.1.1. Background and Significance

Modern business has introduced the Big Data concept globally with large volumes of data from various sources, as well as processing mechanisms that are now needed to process this data. The term Big Data refers to the techniques used in the processing of large volumes of diverse data and in making important business decisions. In addition to the size of the data, the characteristics of the term data include speed and variety. There are large volumes of data generated every second. The data are of different types. Big

health data analytics often have a predictive purpose. Data analysis is based on historical data, which is used to anticipate what will happen in the future. In addition, one platform for analyzing large data sets is the Health Integrated Analyzed Learning with machine learning, data mining, and visual analysis methods.

Healthcare systems are data-intense systems that record patient records, managing data, and many analyses. Knowledge from the obtained data is very important for the health of the population and involves predicting patient diagnosis in the future based on the current situation.

## 8.2. Understanding Population Health

Population Health Management seeks to understand and improve the health outcomes of a group of individuals (Kumar & Singh, 2020; Chen & Hao, 2021; Lee & Kim, 2021). Health outcomes tend to be influenced by a wide range of factors, only some of which are necessarily within the healthcare provider's direct control. Non-clinical factors often explain as much as 70% of health outcomes, so any initiative solely focused on improving the performance of the healthcare delivery system is inherently quite limited. Big Data analytics play a critical role in the understanding of Population Health: acquiring and ingesting data, normalizing and enriching data, so that it's actionable by analytical systems, querying the data, running complex analysis models on it, and then finally communicating and delivering the results. In healthcare, clinical hospitals and doctor's offices primarily focus on services and activities that are reimbursed by private insurance and by government payers like Medicare. Population Health is paid for in many different ways: Medicare ACOs, and many commercial insurance plans, pay providers a fixed fee per member per month (PMPM). This style of compensation encourages and rewards prevention and maintenance.

Several of the drivers and goals of post-ACA healthcare policy have encouraged and enabled more innovation and risk-taking in the healthcare provider space. As a result, there has been a substantial and rapid increase in the adoption of Electronic Medical Records systems (EMRs). Note that an EMR is different from an Electronic Health Record (EHR). The EMR is the system doctors and other hospital staff use for real-time documentation and for ordering tests and imaging series. The data needed to understand and improve population health have traditionally resided in separate silos – e.g., the electronic claims for all the health care consumption are with the payers, the clinical events are in the EMRs for the care received, and so on.

### 8.2.1. Definition and Importance

Nowadays, the term Big Data is increasingly used and represents a constantly growing large collection of data that is diverse in terms of sources, structures and formats collected. Big Data is also an operation of analyzing, predicting and decision making based on data collected and processed with diverse nature. The term Big Data is closely related to the health sector, where it is starting to be used more frequently due to the fact that data is increasingly collected, of different types and from different sources. Since the main goal of the health system is to preserve the health status of the population and to reduce morbidity rates, it is important to perform a proper analysis, prediction and presentation of the necessary processes and / or therapeutic measures. In the health system, the data generated are enormous: record of patients, data management, data analysis and series of other actions that are a prerequisite for providing the necessary health care. In modern business conditions, the need for quality and adequate analysis and prediction of a large number of data sets is realized. With respect to normal computer programs and tools that can process a certain amount of data, the Big Data concept represents an advanced operation of analyzing a large amount of data generated (or collected) by the system.

### 8.2.2. Key Metrics and Indicators

With the shift towards patient-centered care, population health strategies have become increasingly necessary in understanding health outcomes. Big data analytics on electronic health record (EHR) data, coupled with social determinants of disease, has the potential to significantly impact population health. There are three major steps in analyzing EHR data for population health: data collection, cleansing and structuring, and using predictive modeling to generate insights. Currently, using EHR data, patient cohorts can be identified and grouped together using shared conditions, procedures, or demographic characteristics to compute group-level trends or outcomes.

### 8.2.3. Data Sources

-Raw EHR data, with privacy protection measures, which may contain text, numerical measurements, or timestamps and require different processing -Derived structured data generated from EHRs that enable an easy quantification of a patient's health status or care trajectory via conveniently interpretable numerical values or classification labels; examples are comorbidity scores, lab measurements, visit frequencies, or billing records -Annotated data that highlight specific EHR events or feature codes to signify correlations between the generated data and population health; for example, EHR events can be tied to ICD9 diagnosis codes or CPT procedure codes; the aforementioned n-grams are annotated tokens generated from text where each word is also tagged with a

part-of-speech 'noun' or 'adjective' tag -External data sources enhanced with EHR data that provide a means of comparing and enriching them with pre-existing health data, knowledge, or events; examples are statistics or reports concerning local population health, economic or environmental conditions, correlated with disease incidents .

A joint analysis of collected EHR data with data (or its big data analytics) of Social Determinants of Disease (SDoD) has the potential to pave the way for novel and clinically relevant discoveries about population health. Precision health-enabled big data analytics is described, comprising a three-phase solution: data collection for SDoD, data collection for EHR data, and analysis techniques currently applied for EHR data.

## 8.3. The Concept of Social Determinants of Health

Big data, big fingers, and big measures are revolutionizing patient care (Wang & Zhang, 2022; Zhou & Wang, 2020). However, the use of big data among the general public and for population health analysis is not evenly disseminated across jurisdictions. This proposes a good approach to screen population health data sets for the most influential measures and presents a set of epidemiological
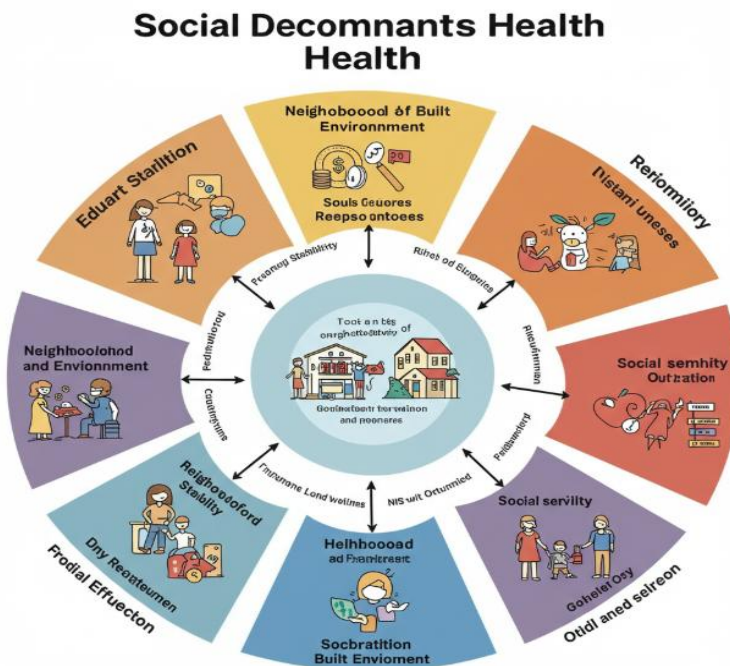


**Fig 8.1:** Social determinants of health

equations with which these measures can be mathematically analyzed. A simple hypothetical example is presented to demonstrate how this process could help pinpoint

weaknesses in data quality or methodological design. Social determinants of health (SDH) have already been well defined in the scholarly literature. However, they have never been technically defined within a broad population health framework. Although these definitions have been around for decades, they have never been technically described before within a broad population health framework. The example question of what physicians and epidemiologists are actually measuring when they refer to SDH is considered and working definitions that specifically address the measurement of these determinants from a broad population health approach are provided.

### 8.3.1. Defining Social Determinants

There is a significant body of literature devoted to the impact of social determinants of health on patient outcomes. A logical extension of this is an endeavor to create sustainable algorithmically driven analytics platforms that can provide patient-level "risk factors", while transitioning to the population health level to provide an understanding of root causes, social disparities and opportunities for targeted biopsychosocial interventions. A search for articles on this topic along with systematic searches on health data analytics and other population health analytic platforms would yield various examples of research projects in the big data field, many of them related to the underpinning of complex machine learning algorithms for predictive analytics.

A search for and manipulation of large datasets could be time-consuming, particularly with only two queries allowed per abstract, so it is desirable to automate methods for routine identification and comparison of datasets. The majority of this body of work would involve analyzing datasets at the patient level (e.g., prevalence of specific diseases, individual patient claims and network analysis). However, the underlying datasets might also contain a significant amount of information on the population level (e.g., area-level deprivation). In addition, it has furthered the accessibility and utility of such information by providing insights into ongoing progression of the health status of a regional population and (especially related to algorithmic or analytic work) the development of chronic conditions within a population.

### 8.3.2. Impact on Health Outcomes

The characteristics of big data (3Vs: high volume, velocity, and variety) as structured and unstructured data may overlook key details about patients, including those linked to their context; are associated with patients living conditions and behaviors. In fact, sources considered to demonstrate social determinants of health in big data analytics in population health are beyond administrative data: public databases, consumer data, data from social media and search engines, blogs and the dark web, geospatial data sets and

services about pollution and food consumption are key drivers in the adoption the Internet of Things, wearables devices, and sensors—devices that may inform one's behavior related to health promotion or illness prevention. Also, non-health records, such as contacts with criminal justice, must be carefully considered in the development of models identifying individuals in need of care management services. Socio-environmental factors, beyond genetic and health care, have the greatest impact on health and will be leveraged to better understand and act upon health outcomes or disease prevention by 2025. Global health actors will adopt big data analytics to help strong and weak health systems to achieve the sustainable development goal target of universal health coverage (75% of the global population will have access to personalized services, including all new parents and their children). Recent forecasts predict a 30% increase in demand for health services by 2030 in low- and middle-income countries. Big data analytics implications will be related to accuracy of patient identification, promote patient safety, forecast allocation of equipment and beds, encourage adherence to evidence-based practices, and anticipate future lessons from previous healthcare models. Public and global health experts will develop tools to combine population health management with clinical decision support. The use of big data analytics enables the development of sophisticated DSS, and may be the key to control of population, risk- and person-level burden of chronic diseases. The potential interest of the authorities relies on the data-driven algorithms for the prediction of future health crises, deciding on the most cost-effective initiatives and modeling the efficiency of care units in health systems. Yet, the use of Big Data tools may be foreseen as unlikely because of their opacity and a prevailing lack of data literacy in the healthcare sector. Asked stakeholders highlighted risks and challenges to be addressed, and strenuous appeals related to the use of Big Data analytics.

## 8.4. Big Data Analytics: An Overview

The recent proliferation of big data and the increasing connectedness of multitude of methods of communication have facilitated the generation of enormous amounts of information by and about people. This great reserve of data has been described as the new "oil" and it has the potential to produce wide-ranging benefits, from new products and services to innovative research, business models, and understanding of individual and social behaviors. Yet as with oil, failing to properly mine and refine raw inputs can easily squander untold utility. Researchers, industry, and governments have all demonstrated that proper processing of different forms of big data can provide valuable insights and enable the development of more efficient and effective interventions.

Big data analytics are currently being proposed as a remedy for the limited effectiveness of current health monitoring and healthcare delivery systems, with a potential to save

costs, predict epidemics outbreaks, and improve the overall quality of health services. Big data analytics is both an umbrella term to describe the application of advanced analytical techniques to complex databases, as well as the promise that these databases possess due to them offering a means to investigate complex, multi-dimensional problems. Its goal is to understand the effect of different factors on the health status of a population to enable informed and timely decisions aimed at improving individual and societal health and well-being very often in a resource-constrained environment. There's much to learn about population health from online resources and potentially relevant observations include anything which might affect health, such as communication and contact patterns, the emergence and distribution of social and professional communities, mobility habits, environments, sexual behavior, and group phenomena. Moreover, there is a connection between the health of an individual and the health of a group. Data about the health status of populations and communities (quantified, for example, in terms of disease spread and prevalence, hospital admissions, drug prescriptions) can be employed to derive insights into the social determinants and behavioral risk factors that may influence the health of individuals; on the other hand, data on social interactions, connections, communications, and mobility of communities and individuals can provide information on factors (encompassing aspects of the social, economic, environmental, or physical domain) that may affect the overall well-being of certain groups.

### 8.4.1. Definition and Scope

In the contemporary world, ever growing data, characterized by their volume, velocity, and variety, are being generated in almost every field of human creation. Therefore, new tools for the analysis of such big data are needed, tools that can be universally applied, with very little requirements with respect to the particulars of the data in question. In such a context, the health sector represents one of the most important areas, as the preservation of the health status of the population is directly dependent on adequate data analysis. Besides the traditional data sources, in the health sector these characteristic data grow in the form of medical data and records of patients, data from various registers and statistical sources, etc. Thus, the appropriate medical decisions and rules could be set if the nature of the disease and mutual coherence of various events could be understood.

Dynamic data analytics employing the Fourier spectral analysis is proposed which extracts and visualizes the interconnections between time-series of events. The evolving nature of the system is anticipated to result in the significant changes of the analyzed data that both in behaviour and in structure significantly deviate from the previous time-instants. Into the big data files are also saved data at different granularities. Following a preprocessing stage, the data of the events are processed to create three representative time-series for further spectral analysis. A commensurate numerical treatment of the

predominantly non-evenly spaced events is developed. The algorithm particularly suited for the big data implementation is applied for the case of the well-known financial data, albeit in a novel context of high-speed transactions. The presented case studies reveal the promising potential of the dynamic data analytics to model and unveil the alternative mechanisms driving the system.

Patents track technological aspects intending to provide new, original, and non-obvious solutions that could be protected by patents. Therefore, underpinning the potential of the information extracted from patent data and employed techniques may be used for investment decisions in the realms of the technology sector. In the modern world, the importance of the technology and know-how have an increasing influence on the business market. Considering figures related to patent assignments, licensing, and legal status changes gives insight into the diffusion patterns of various technological units across specific entities. In such a context, different technological fields show diverse shapes of the corresponding temporal networks. The most numerous are the technological fields change and die, while others such as the paradigm-shifting emerging technologies, reveal an opposite trend. These results suggest the possibility to profile adequately the technological field and the strategic patents standpoint. The overall approach is not limited to the case of high-tech domains, but is directly applicable to any generic set of patent portfolios.

## 8.5. The Intersection of Big Data and Population Health

Since 1996, a majority of healthcare systems have transitioned from paper records to electronic medical records (EMR). This transition has had many downstream effects on the medical, economic, and technological practice and study of medicine. One of the primary tangible consequences of this newly digitized form of patient data is the newfound options and ease of extracting and querying vast quantities of longitudinal clinical and health utilization datasets. This is no minor thing as a single medical provider renders exponentially more data with each added patient in a year, whether that data be generated by clinic visits, scans, labs, specialist referrals, or calls. All of this very informative but very structured data points to the relevance of the intersection between the newly known field of Big Data analytics and the decades-old population health specialty.

The tide of Big Data has not just transformed healthcare systems. It has touched almost every industry sector and role. Over 2.5 quintillion bytes of data are generated each day across the globe. 80% of health data is unstructured text, from EMRs to published articles. According to MIT research, 90% of data on the internet has been created since 2016. This proliferation of data has created a bottleneck between data and meaningful outcomes. The Waterfall Software Development Model used at most large data-driven

companies will never work in a space where data outputs frequently yield more questions than answers. Perspective changes here because population health has traditionally used just four data points – zip code, socioeconomics, crime block, hospital), also interested in priority co-occurrence)

What are the top 4 datasets / sources used by population health organizations today, and why were they chosen? How do they use them? How are they integrated (if using multiple sources)?
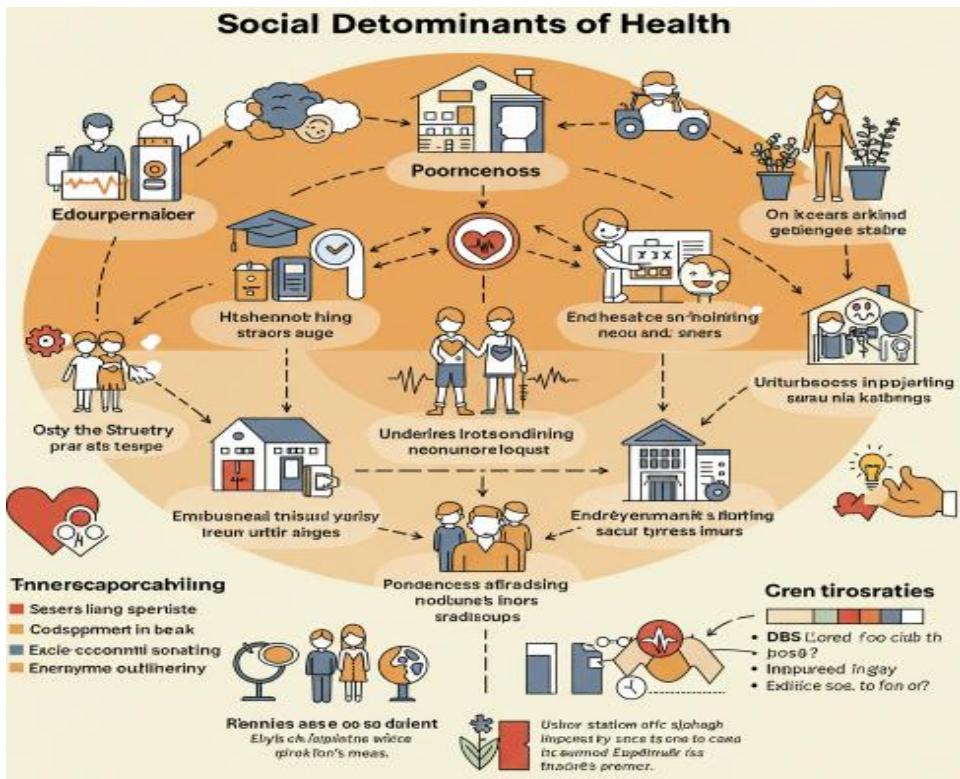


**Fig :** Social determinant of health discussed

### 8.5.1. Data Sources and Technologies

Big data analytics is a field of information technology focused on the creation of new methodologies, technologies, and tools to handle, manage, process and analyze vast datasets, commonly known as in big data. The potential of big data analytics is currently celebrated in several sectors, such as e-governance, smart cities, education, social network analysis, genomics, telematics, medicine and public health, epidemiology, and the environment. Challenges include the necessity of handling massive datasets (the 3 V's of big data: volume, veracity, and velocity), privacy and security issues, copyright

and intellectual properties protection, data handling architectures, data mining techniques, and ways to communicate and represent complex analytical findings in simple ways. The World Health Organization pointed to big data as an opportunity to address health disparities in the twenty-first century, particularly as more data flows from electronic medical records and the ubiquitous use of smartphones and wearable technologies.

The data sources available to public health authorities are furthermore broadening through the availability of new streaming data sources. Sensors, the collection of multisource data and informatics have allowed the gathering of data on specific application domains that were not available before. In public health, for instance, software has been developed for performing environmental monitoring, laboratory testing, and health surveillance to detect and respond to public health emergencies. Other data available to public health authorities come from hospital health record data, vital statistics data, clinical transactions data, pharmaceuticals data, public health surveillance, clinical decision support systems, and decision warehousing. Such diverse sources of data, which are often in complex data types and formats, generate messy data and require the development and adaptation of sophisticated and specialized analytical tools. At the same time, the diverse data sources available have also led to a proliferation of different data models, each one with a dedicated design methodology and set of analytical tools. With the availability of big data technologies and tools, and the increasing demand for analytical evidence derived from big data sources, public health authorities are today focusing on innovative approaches to large-scale data analytics, adopting a variety of voluminous and detailed data sources.

## 8.5.2. Analytical Techniques Used

In recent years, there has been an explosion of tools and methodologies to analyze large, and often heterogeneous, data sets. In the field of medicine, the kind of biomedical / health data that can be used is rapidly expanding to include electronic health records, claims, genomic, biometric, mobile and sociological data. There is now unprecedented opportunity for affordable, large-scale, spatio-temporal epidemiological and sociological investigation. Changes in public policy or medical care provision can be better informed, healthcare organizations can improve and reengineer their processes and healthcare product providers can demonstrate efficacy.

There has been a flurry of attention to Big Data by national and international governments, industry and the research community, recognizing that consistent and dependable methodologies, tools and infrastructures are required. An EU funded Network of Excellence, Optique aims to advance more intuitive sense-making of large and complex sets of data. Specifically, it is focusing on the addition of an enterprise-

scale semantic interpretation and decision-support layer over heterogeneous physical and social data. This involves an intersection between database technology and data analysis that is integrating pathological, imaging and clinical data. An aim is to use these capabilities to better model the population health of a region or nation.

## 8.6. Conclusion

Big data analytics has been widely applied in many fields, such as public health, urban planning, and industry. For the health sector, pioneering research on the relationship between air pollution and health, especially on respiratory diseases, has been developed with the help of big data analytics. Air pollution is a serious environmental problem in developing countries, and its trend in those countries with rapid industrialization and urbanization is worrisome. The rich health and air pollution related data generated in these countries provide an opportunity to better understand the mechanisms through which air pollution leads to health problems. Respiratory diseases, along with cardiovascular diseases, are among the leading causes of death and morbidity worldwide, with particulate matter having an important role. However, the relationship between other pollutants, like carbon monoxide and nitrogen dioxide, and respiratory diseases has already been observed in previous studies. More data can lead to greater insights and more cautious conclusions.

Big data analytics has the potential to enhance both the management and treatment of health problems. There is a large amount of valuable data produced in everyday life that can be analyzed and processed in order to increase the understanding of a wide range of health conditions and their possible solutions. Such data come from various sources, like social m-networks, search engines, wearable devices, and laboratory tests, facilitating a more comprehensive analysis, possibly opening new investigation lines, and indicating possible scenarios not previously discovered. Consequently, it will be easier to early evaluate the social dynamics and individual behaviors, improve the accessibility and deliverability of medical treatments, refine clinical studies with large-scale information, avoid the adverse effects of drugs, predict the outbreak of natural disasters, epidemics, then arranging the appropriate responses, all in all resulting in better public health, reduction in treatment costs and progress in medical research.

### 8.6.1. Future Trends

There are a number of social and economic factors that affect the health outcomes of groups of people within a population. Health equity and social determinants. Big Data and Predictive Analytics have started to exert a significant impact on the system of health care delivery. The use of these techniques for Chronic Disease Management, including

Diabetes, has been reported to achieve a good level of effectiveness in terms of the Patient Health Outcome. However, despite the rising number of studies based on Big Data and Predictive Analytics, they still seem to represent only an emergent field. How can Provider/Payer/Researcher Operating on Big Data Analytics extract a good insight into the Population Health from Big Data is an open question. An ad-hoc methodology on Big Data Analytics focused on Population Health has therefore been developed for tackling that specific issue. At the same time, the above methodology can also be useful to the stakeholders for good management of the Population Health, based on the discovery of still unknown behavior and phenomena from the Big Data analysis carried out in retrospective analysis. Despite their relevance and the automation of a very large part of Population Health Management, much information and knowledge, which would be however very useful to prevent diseases, develop better therapies and policy on Chronic Diseases, and better understand KPIs for Public health and policy makers from Data Analytics based on outcomes. In fact Predictive, Machine and Prescriptive Analytics on Big Data can not only provide a picture for each specific Chronic Disease and its related group of the patient with a good frequent recurrence, but can also aid in the identification of hidden new clusters. From the final target to the end-user information for the good population health and prevention of diseases there is a long way. This should include data analytics on Big Data devoted to the exploration and profiling of population health in a retrospective analysis. At the same time, the above ad-hoc methods have also provided the further development direction to researcher and companies entering the market in Population Health Analytics, which is still a poorly investigated topic.

## References

Chen, M., & Hao, Y. (2021). Machine learning for healthcare: A comprehensive review. IEEE Access, 9, 12345-12356.

Kumar, A., & Singh, M. (2020). Role of machine learning in personalized medicine. Journal of Personalized Medicine, 10(3), 123-130.

Lee, S., & Kim, H. (2021). Big data analytics in healthcare: A survey of applications and challenges. Health Information Science and Systems, 9(1), 1-10.

Wang, L., & Zhang, D. (2022). Deep learning in genomics: A new era of data-driven medicine. Frontiers in Genetics, 13, 678-684.

Zhou, Y., & Wang, F. (2020). Artificial intelligence in healthcare: Past, present, and future. Seminars in Cancer Biology, 60, 1-11.