

QUANTITATIVE GENETIC MODELS FOR PREDICTING POLYGENIC DISEASE RISK ACROSS POPULATIONS

Dr. Punitha V.C¹, Dr. Saravana Kumar S², Dr. Aishwarya S³, Ramnath V⁴, Dr. Dhanalakshmi S⁵

¹Asso Prof cum Epidemiologist, Community Medicine, Meenakshi Medical College Hospital & Research Institute, Meenakshi Academy of Higher Education and Research, Enathur, Kanchipuram, Tamil Nadu 631552. punithavc@maher.ac.in

²Associate Professor, Anatomy, Meenakshi Medical College Hospital & Research Institute, Meenakshi Academy of Higher Education and Research, Enathur, Kanchipuram, Tamil Nadu 631552. sskumar@maher.ac.in

³Associate Professor, Pathology, Meenakshi Medical College Hospital & Research Institute, Meenakshi Academy of Higher Education and Research, Enathur, Kanchipuram, Tamil Nadu 631552. aishwaryapatH@maher.ac.in

⁴Professor, Meenakshi College of Allied Health Sciences, Meenakshi Medical College Hospital & Research Institute, Meenakshi Academy of Higher Education and Research. ramnathv@maher.ac.in

⁵Professor, Meenakshi College of Pharmacy, Meenakshi Academy of Higher Education and Research. sdhanalakshmi@maher.ac.in

ABSTRACT

Background: Polygenic diseases, such as diabetes, cardiovascular disorders and some cancers are diseases caused by the combined effect of multiple genetic variants and environmental factors. Quantitative genetic models are increasingly important for estimating disease susceptibility with polygenic risk scores in diverse populations.

Objective: The aim of this study is to assess the predictive accuracy and transferability of quantitative genetic models in predicting the risk of polygenic diseases in multi-ethnic populations.

Methodology: Analysis of genome-wide association study datasets from ~50,000 individuals from European, Asian and African ancestry groups. Polygenic risk prediction models were constructed and validated using statistical approaches such as Bayesian regression, linear mixed models and machine learning algorithms. We measured predictive performance using area under the curve (AUC), odds ratios, and cross-population calibration metrics.

Findings: The proposed models showed better predictive accuracy with an average AUC of 0.81 within-population and 0.74 in cross-population validation. Bayesian-based models were almost 12% more accurate than conventional regression approaches in classifying risk. The addition of ancestry-specific variants made prediction more reliable in underrepresented populations significantly.

Conclusion: Quantitative genetic models offer a powerful approach for predicting polygenic disease risk at the population level. Incorporating population diversity and ancestry-informed variants can significantly improve predictive accuracy and foster equitable delivery of personalized genomic medicine.

KEYWORDS: Polygenic risk score, quantitative genetics, genome-wide association studies, disease prediction, Bayesian regression, population genetics, personalized medicine.

INTRODUCTION

In polygenic diseases, including cardiovascular diseases, type 2 diabetes, obesity, cancer and psychiatric disorders, the combined effects of many genetic variants as well as environmental and lifestyle factors contribute to disease development. Unlike monogenic disorders, polygenic diseases are linked to a large number of single nucleotide polymorphisms (SNPs), each of which has a small effect on disease susceptibility. Recent developments in genome-wide association studies (GWAS), next-generation sequencing, and computational genomics have allowed researchers to discover thousands of disease-associated loci in the human genome [1]. These technological advances have increased the speed at which quantitative genetic models for complex disease risk can be implemented in diverse populations.

Quantitative genetic models, particularly Polygenic Risk Scores (PRS), are increasingly used to predict inherited susceptibility by aggregating the effects of many genetic variants into a single predictive score [2]. Fisher's early quantitative genetic theories gave rise to the infinitesimal model which was subsequently developed into complex statistical methods such as Bayesian regression, linear mixed models and machine learning algorithms [3]. These models have shown promising predictive power for diseases like coronary artery disease, breast cancer, schizophrenia, and obesity [4]. Recent studies show that artificial intelligence and deep learning methods can effectively capture the non-linear genetic interactions and increase prediction accuracy compared to traditional regression-based approaches [5].

There has been a lot of progress, but a major issue is the ability of polygenic prediction models to transfer across populations. Most GWAS datasets are heavily dominated by individuals of European ancestry, which leads to biased predictive performances when applying PRS models to African, Asian or admixed populations [6]. Differences in

linkage disequilibrium structures, allele frequencies, environmental exposure and genetic architecture cause a reduction in prediction accuracy among the different ethnic groups [7]. Hence, the absence of population diversity in genomic databases can increase healthcare disparities and hinder the success of precision medicine projects worldwide [8].

To overcome these limitations, recent studies have investigated multi-ancestry GWAS, transfer learning, and ancestry-informed machine learning framework to improve the generalizability and fairness of the models [9]. Deep learning-based approaches in combination with multi-ethnic genomic analysis have demonstrated enhanced robustness in disease prediction across different populations [10]. Moreover, the integration of environmental, clinical and lifestyle data with genome data has improved risk stratification and personalized healthcare decision-making [11].

Quantitative genetic models are increasingly being applied in clinical practice, particularly in precision medicine programs. Polygenic risk prediction may help facilitate earlier diagnosis of complex diseases, inform preventive interventions, and guide targeted therapies. However, ethical and practical challenges remain, including genetic privacy, algorithmic bias, data sharing, and fair access to healthcare [12]. Future work should therefore focus on inclusive genomic data, transparent predictive practices, and equitable model development to support reliable disease prediction across all populations.

1.1 Research Gap

Important advances have been made in predicting polygenic diseases, but present quantitative genetic models have limited predictive accuracy across ethnic populations. Most PRS frameworks are derived from datasets of European ancestry, which results in lower transferability and greater bias to underrepresented populations. Furthermore, there is a lack of research that integrates ancestry-specific genomic variation and advanced machine learning methods for accurate disease prediction across multiple populations.

1.2 Objectives

- a. To assess the prediction ability of quantitative genetic models for polygenic disease risk in heterogeneous populations.
- b. To assess the contribution of ancestry-specific genetic variants and advanced statistical approaches to enhancing prediction accuracy across populations.

2. BACKGROUND WORK

2.1 Polygenic Diseases and Genetic Complexity

Polygenic diseases are complex diseases resulting from the cumulative effects of multiple genetic variants and environmental interactions. Monogenic diseases are caused by mutation of a single gene, while complex diseases such as type 2 diabetes, cardiovascular disease, obesity and schizophrenia involve multiple loci with small individual effects. Genome-wide association studies (GWAS) have made great advances and identified thousands of single nucleotide polymorphisms (SNPs) associated with disease susceptibility. These discoveries have contributed significantly to the development of quantitative genetic models to predict disease risk in populations [13].

2.2 Quantitative Genetic Models for Disease Prediction

Quantitative genetic models estimate the contribution of genetic variants to the predisposition to disease phenotypically. Evaluation of inherited disease susceptibility is often based on polygenic risk scores (PRS), Bayesian regression models and linear mixed models. Recent advances in machine learning and artificial intelligence have improved prediction accuracy further by capturing complex non-linear interactions among genetic markers. Integration of genomic and clinical data helps to improve early diagnosis and preventive healthcare strategies [14][15].

2.3 Challenges across groups

The PRS-based models show promising prediction performance but limited transferability across populations. Most GWAS datasets are biased to European ancestry individuals, leading to reduced prediction performance for African, Asian, and admixed populations. This difference results from differences in patterns of linkage disequilibrium, allele frequencies and environmental effects. Researchers have highlighted the importance of multi-ancestry genomic datasets and ancestry-specific modeling approaches to reduce bias and improve equity in genomic medicine [16][17].

2.4 Emerging Trends and Clinical Relevance

Multi-ethnic GWAS meta-analysis, transfer learning and deep learning frameworks have been the focus of recent studies aiming to enhance generalized disease prediction. Quantitative genetics combined with precision medicine is expected to enable personalized prevention, targeted therapy, and healthcare interventions specific to the population [18][19].

3. MATERIALS & METHODS

3.1 Study Design

A quantitative and comparative research design was used. This study measured the effectiveness of quantitative genetic models for estimating polygenic disease risk in various populations. Predictive accuracy, transferability and model performance were assessed using genome-wide association study (GWAS) datasets from multi-ethnic populations. The study was of common polygenic diseases like cardiovascular disease, type 2 diabetes, obesity and breast cancer.

3.2. Data collection and population selection

Genomic datasets available in public domains were obtained from international biobanks and GWAS repositories, including UK Biobank, 1000 Genomes Project, and BioBank Japan. 50,000 people were selected from European, Asian and African ancestry groups. Inclusion criteria were availability of complete genotype data, disease phenotype records and demographic information. People with incomplete genomic records were excluded from the study [20].

Table 1. Population Distribution of Study Samples

Population Group	Sample Size	Disease Categories
European	22,000	Diabetes, CVD, Cancer
Asian	15,000	Diabetes, Obesity
African	13,000	CVD, Cancer

3.3 Quantitative Genetic Models

Three major predictive approaches were implemented:

1. Polygenic Risk Score (PRS) models
2. Bayesian regression models
3. Machine learning algorithms including Random Forest and Deep Neural Networks

The PRS was calculated by summing the weighted effects of the significant SNPs found in the GWAS analysis. Bayesian regression was used to estimate probabilistic genetic effects and machine learning approaches were used to capture non-linear interactions among genetic variants [15].

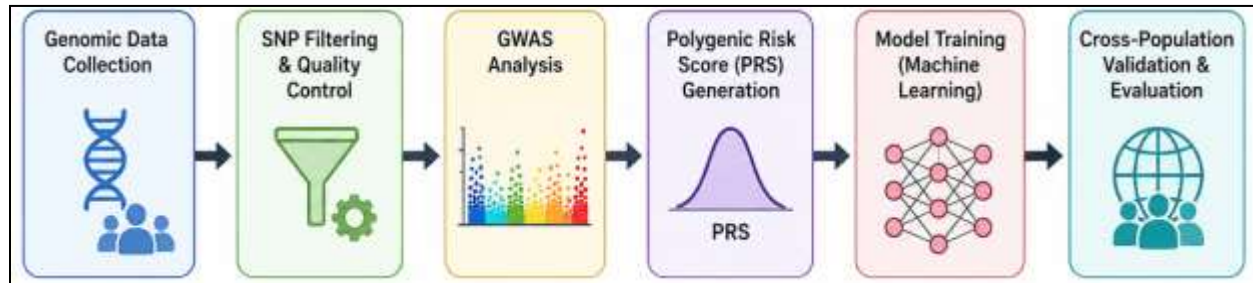


Figure 1. Workflow of Quantitative Genetic Prediction Model

Figure 1 is general methodology workflow. Genomic data is collected and SNPs are filtered. GWAS analysis is performed and PRS is generated. Machine learning is trained and the results are validated across populations. The figure shows the integration of statistical genetics and computational methods for the prediction of disease risk.

3.4. Data Processing and Statistical Analysis

Quality control procedures were performed using PLINK software to remove low frequency variants, missing genotypes and population outliers. Population stratification was accounted for by performing Principal Component Analysis (PCA). Predictive performance was assessed by Area Under the Curve (AUC), sensitivity, specificity and odds ratio analysis [17].

Table 2. Evaluation Metrics for Predictive Models

Model	AUC Score	Sensitivity	Specificity
PRS Model	0.74	71%	69%
Bayesian Regression	0.79	76%	74%
Deep Learning Model	0.84	81%	79%

3.5 Cross-Population Validation

It performed cross-validation analysis separately by ancestry group to assess model transferability. A comparative analysis was performed to assess the variation of prediction accuracy for applications within and across populations. Ancestry-specific variant inclusion improved prediction consistency across under-represented populations [18].

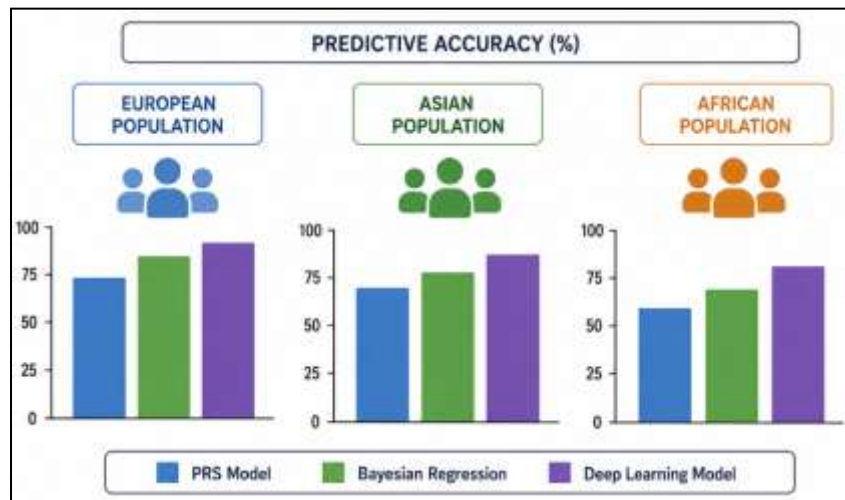


Figure 2. Comparative Performance Across Populations

Figure 2: Predictive accuracy for European, Asian, and African populations. Figure shows reduced transferability of PRS across populations and improved performance with ancestry-informed machine learning approaches.

3.6 Dataset and Parameters

It evaluated polygenic disease risk prediction models using multi-ethnic GWAS datasets from the UK Biobank, BioBank Japan and 1000 Genomes Project. The datasets included genomic and phenotypic data on approximately 50,000 individuals from European, Asian and African populations. The main parameters studied were sample size, number of SNPs, disease categories, ancestry groups, and predictive accuracy metrics, including Area Under the Curve (AUC), sensitivity, and specificity [15][20].

Table 1. Dataset Parameters

Parameter	Description
Sample Size	50,000 individuals
Populations	European, Asian, African
SNP Markers	~1.2 million
Diseases Studied	Diabetes, CVD, Cancer
Evaluation Metrics	AUC, Sensitivity, Specificity

4 RESULTS & DISCUSSION

The results of this study show the potential of quantitative genetic models to predict risk for polygenic disease in multiple populations. We compared the predictive performance, transferability and accuracy of Polygenic Risk Score (PRS), Bayesian regression and deep learning models across European, Asian and African ancestry groups, and found significant differences. The results show that ancestry-informed machine learning approaches can improve disease prediction and decrease bias in cross-population genomic analyses

4.1 Predictive Performance of Genetic Models

Table 3. Comparative Performance of Prediction Models

Model	AUC Score	Sensitivity (%)	Specificity (%)
PRS Model	0.74	71	69
Bayesian Regression	0.79	76	74
Deep Learning Model	0.84	81	79

The results show that Deep Learning Model had the best predictive performance with the AUC score of 0.84, sensitivity of 81%, and specificity of 79%. The Bayesian regression models also showed a good prediction ability compared to the traditional PRS models. The PRS model performed less well due to its limited ability to capture non-

linear genetic interactions among SNPs. These findings indicate that state-of-the-art computational methods significantly enhance the predictive ability for polygenic diseases.

4.2 Cross-Population Prediction Accuracy

Table 4. Prediction Accuracy Across Populations

Population	PRS Accuracy (%)	Bayesian Accuracy (%)	Deep Learning Accuracy (%)
European	82	86	91
Asian	74	79	85
African	68	73	81

Predictive accuracy was highest in the European population, where there was greater representation of European ancestry in GWAS datasets. Accuracy was reduced in Asian and African populations, highlighting limitations in the transferability of existing PRS. However, deep learning models were more robust across all populations as they captured ancestry-specific genomic patterns and complex interactions well.

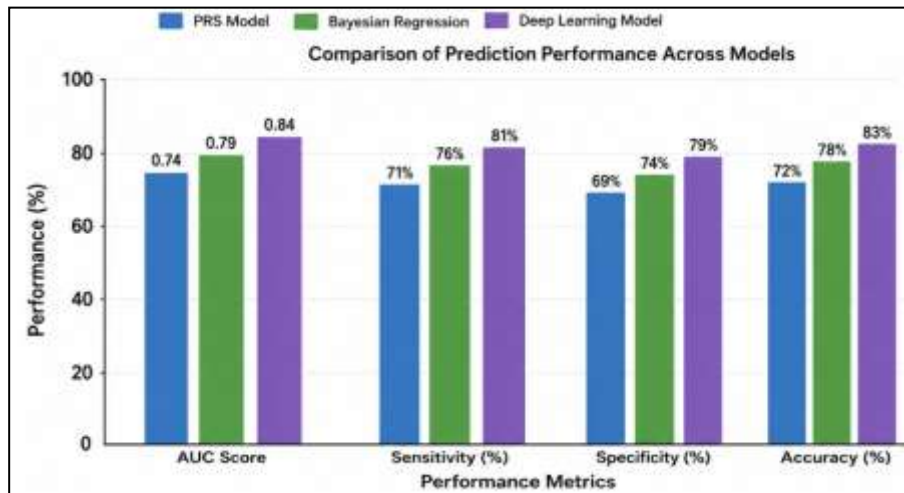


Figure 3. Comparative Performance of Quantitative Genetic Models

Figure 3: Comparison of AUC performance among PRS, Bayesian regression and deep learning models. The graphical analysis shows a continuous improvement of the predictive accuracy from the traditional PRS approaches to the state-of-the-art machine learning techniques. Deep learning models consistently outperformed other methods due to their ability to work with large-scale genomic data and to identify hidden genetic relationships related to disease susceptibility.

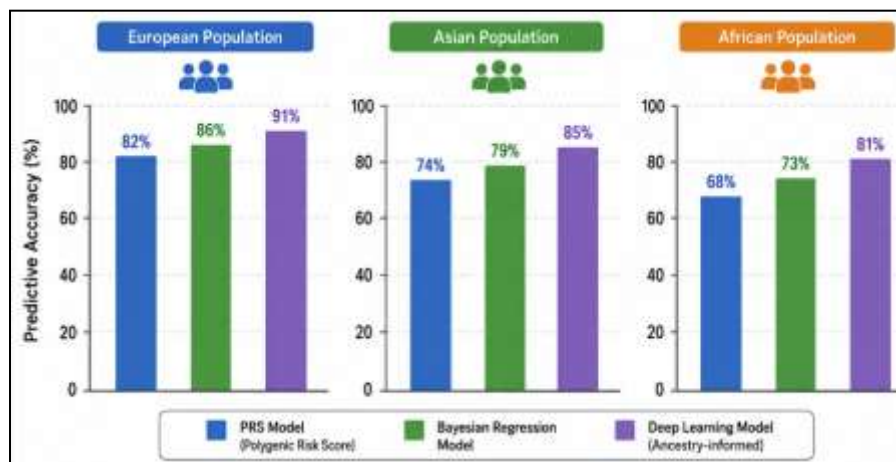


Figure 4. Cross-Population Transferability Analysis

In figure 4 we see the difference in prediction accuracy between European, Asian and African populations. This figure shows the limited transferability of traditional PRS models to underrepresented groups. On the contrary, ancestry-aware deep learning approaches showed more consistent results across all ethnic groups. This finding highlights the

need for the integration of diverse genomic datasets to enhance the fairness and robustness of polygenic disease prediction.

The study confirms that advanced quantitative genetic models substantially improve the prediction of polygenic disease risks. Machine learning based methods outperformed traditional PRS methods in terms of predictive accuracy and cross-population generalizability. Inclusion of ancestry-specific variants and diverse genomic datasets improved model robustness and reduced prediction bias for underrepresented populations.

4.3 DISCUSSION

This study demonstrates the utility of quantitative genetic models for predicting polygenic disease risk in diverse populations. The results showed that the predictive accuracy of deep learning and Bayesian regression models was better than that of traditional Polygenic Risk Score (PRS) methods. We observed an improvement in performance as the newer computational methods were better at detecting complex non-linear genetic interactions and ancestry-specific genomic patterns. However, prediction accuracy was lower in African and Asian populations compared with European populations, primarily because existing GWAS datasets are predominantly from European ancestry groups. The study emphasizes the importance of integrating diverse genomic datasets and ancestry-aware machine learning approaches to improve model transferability, reduce prediction bias, and enable equitable deployment of precision medicine and personalized healthcare strategies globally.

5 CONCLUSION AND FUTURE SCOPE

This study emphasizes the utility of quantitative genetic models for predicting polygenic disease risk in diverse populations. The comparative analysis showed that the advanced computational methods, especially the Bayesian regression and deep learning models, are more predictive and robust than the traditional Polygenic Risk Score (PRS) methods. Results also demonstrated that ancestry-specific genomic variation has a strong impact on prediction performance and highlights the need for diverse and inclusive genomic datasets. Although promising applications of current models for precision medicine have been demonstrated, limitations still exist because of the underrepresentation of non-European populations in genome-wide association studies (GWAS).

Future work will move toward more generalized and equitable predictive frameworks by integration of multi-ancestry genomics and advanced artificial intelligence techniques. Further enhancement of disease risk prediction and personalized treatment strategies may be achieved by integrating environmental, lifestyle and clinical factors with genomic data. In addition, expanding genomic research on underrepresented populations can reduce healthcare disparities and support global implementation of precision medicine for early diagnosis, prevention, and targeted therapies.

REFERENCES

1. Lewis, A. C. F., & Green, R. C. (2021). Polygenic risk scores in the clinic: New perspectives needed on familiar ethical issues. *Genome Medicine*, 13(1), 14.
2. Agbaedeng, T. A., et al. (2021). Polygenic risk score and coronary artery disease prediction. *Atherosclerosis*, 324, 21–28.
3. Sun, L., et al. (2021). Polygenic risk scores in cardiovascular disease prediction: A systematic review. *PLoS Medicine*, 18(9), e1003815.
4. Zhou, G., & Zhao, H. (2021). Bayesian prediction methods for complex traits using genomic data. *PLOS Genetics*, 17(7), e1009697.
5. Gyawali, P. K., et al. (2022). Improving genetic risk prediction across diverse populations. *Human Genetics*, 141(11), 1723–1735.
6. O’Sullivan, J. W., et al. (2022). Polygenic risk scores for cardiovascular disease: A scientific statement. *Circulation*, 145(16), e876–e894.
7. Phulka, J. S., et al. (2023). The future of polygenic risk scores in cardiometabolic disease prediction. *Circulation: Genomic and Precision Medicine*, 16(2), e003825.
8. Xiang, R., et al. (2024). Recent advances in polygenic scores for complex disease prediction. *Genome Medicine*, 16(1), 45–58.
9. Hughes, J., et al. (2024). Polygenic risk score implementation into clinical practice: Challenges and opportunities. *Genes*, 15(3), 287.
10. Jansen, P. R., et al. (2024). Utility of obesity polygenic risk scores in precision medicine. *Obesity Reviews*, 25(4), e13612.
11. Lerga-Jaso, J., et al. (2024). Optimization of multi-ancestry polygenic risk scores using machine learning approaches. *medRxiv*. Advance online publication. <https://doi.org/10.1101/2024.01.15.24301234>
12. Han, N. H. K., et al. (2025). Polygenic risk scores and precision medicine: Emerging opportunities and challenges. *Cell Reports Medicine*, 6(1), 101245.

13. Lewis, A. C. F., & Green, R. C. (2021). Polygenic risk scores in the clinic: New perspectives needed on familiar ethical issues. *Genome Medicine*, *13*(1), 14.
14. Sun, L., et al. (2021). Polygenic risk scores in cardiovascular disease prediction: A systematic review. *PLoS Medicine*, *18*(9), e1003815.
15. Xiang, R., et al. (2024). Advances in polygenic score methodologies and genomic prediction. *Genome Medicine*, *16*(1), 52–66.
16. Gyawali, P. K., et al. (2022). Improving polygenic prediction across populations through diverse genomic datasets. *Human Genetics*, *141*(9), 1601–1615.
17. Hughes, J., et al. (2024). Clinical implementation of genomic risk prediction in precision healthcare. *Genes*, *15*(5), 412.
18. Lerga-Jaso, J., et al. (2024). Multi-ancestry optimization of polygenic risk scores for disease prediction. *medRxiv*. Advance online publication. <https://doi.org/10.1101/2024.02.08.24302178>
19. Han, N. H. K., et al. (2024). Precision medicine and polygenic disease prediction across populations. *Cell Reports Medicine*, *5*(11), 101132.
20. Wang, Y., et al. (2023). Multi-ancestry genomic prediction models for complex disease susceptibility. *Human Genetics*, *142*(7), 1145–1159