

Breedon, JR, Marshall, CR, Giovannoni, G, van Heel, DA, Akhtar, S, Anwar, M, Arciero, E, Asgar, O, Ashraf, S, Breen, G, Chung, R, Curtis, CJ, Chaudhary, S, Chowdhury, M, Colligan, G, Deloukas, P, Durham, C, Durrani, F, Eto, F, Finer, S, Garcia, AA, Griffiths, C, Harvey, J, Heng, T, Huang, QQ, Hurles, M, Hunt, KA, Hussain, S, Islam, K, Jacobs, BM, Khan, A, Khan, A, Lavery, C, Lee, SH, Lerner, R, MacArthur, D, Malawsky, D, Martin, H, Mason, D, Mazid, MB, McDermott, J, McSweeney, S, Miah, S, Munir, S, Newman, B, Owor, E, Qureshi, A, Rahman, S, Safa, N, Solly, J, Tahmasebi, F, Trembath, RC, Tricker, K, Uddin, N, van Heel, DA, Winckley, C, Wright, J, Dobson, R and Jacobs, BM

Polygenic risk score prediction of multiple sclerosis in individuals of South Asian ancestry

<http://researchonline.ljmu.ac.uk/id/eprint/19008/>

Article

Citation (please note it is advisable to refer to the publisher's version if you intend to cite from this work)

Breedon, JR, Marshall, CR, Giovannoni, G, van Heel, DA, Akhtar, S, Anwar, M, Arciero, E, Asgar, O, Ashraf, S, Breen, G, Chung, R, Curtis, CJ, Chaudhary, S, Chowdhury, M, Colligan, G, Deloukas, P, Durham, C, Durrani, F, Eto, F, Finer, S, Garcia, AA, Griffiths, C, Harvey, J, Heng, T, Huang, QQ.

LJMU has developed **LJMU Research Online** for users to access the research output of the University more effectively. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in LJMU Research Online to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

The version presented here may differ from the published version or from the version of the record. Please see the repository URL above for details on accessing the published version and note that access may require a subscription.

For more information please contact researchonline@ljmu.ac.uk

BRAIN COMMUNICATIONS

Polygenic risk score prediction of multiple sclerosis in individuals of South Asian ancestry

Joshua R. Breedon,¹ Charles R. Marshall,^{1,2} Gavin Giovannoni,^{1,2,3} David A. van Heel,³ Genes & Health Research Team Ruth Dobson^{1,2*} and Benjamin M. Jacobs^{1,2,*}

* These authors contributed equally to this work.

Polygenic risk scores aggregate an individual's burden of risk alleles to estimate the overall genetic risk for a specific trait or disease. Polygenic risk scores derived from genome-wide association studies of European populations perform poorly for other ancestral groups. Given the potential for future clinical utility, underperformance of polygenic risk scores in South Asian populations has the potential to reinforce health inequalities. To determine whether European-derived polygenic risk scores underperform at multiple sclerosis prediction in a South Asian-ancestry population compared with a European-ancestry cohort, we used data from two longitudinal genetic cohort studies: Genes & Health (2015–present), a study of ~50 000 British–Bangladeshi and British–Pakistani individuals, and UK Biobank (2006–present), which is comprised of ~500 000 predominantly White British individuals. We compared individuals with and without multiple sclerosis in both studies (Genes & Health: $N_{\text{Cases}} = 42$, $N_{\text{Control}} = 40\,490$; UK Biobank: $N_{\text{Cases}} = 2091$, $N_{\text{Control}} = 374\,866$). Polygenic risk scores were calculated using clumping and thresholding with risk allele effect sizes obtained from the largest multiple sclerosis genome-wide association study to date. Scores were calculated with and without the major histocompatibility complex region, the most influential locus in determining multiple sclerosis risk. Polygenic risk score prediction was evaluated using Nagelkerke's pseudo- R^2 metric adjusted for case ascertainment, age, sex and the first four genetic principal components. We found that, as expected, European-derived polygenic risk scores perform poorly in the Genes & Health cohort, explaining 1.1% (including the major histocompatibility complex) and 1.5% (excluding the major histocompatibility complex) of disease risk. In contrast, multiple sclerosis polygenic risk scores explained 4.8% (including the major histocompatibility complex) and 2.8% (excluding the major histocompatibility complex) of disease risk in European-ancestry UK Biobank participants. These findings suggest that polygenic risk score prediction of multiple sclerosis based on European genome-wide association study results is less accurate in a South Asian population. Genetic studies of ancestrally diverse populations are required to ensure that polygenic risk scores can be useful across ancestries.

- 1 Preventive Neurology Unit, Wolfson Institute of Population Health, Queen Mary University of London, London EC1M 6BQ, UK
- 2 Department of Neurology, Royal London Hospital, London E1 1FR, UK
- 3 Blizard Institute, Queen Mary University of London, London E1 2AT, UK

Correspondence to: Dr. Ruth Dobson
Preventive Neurology Unit
Wolfson Institute of Population Health, Charterhouse Square
London EC1M 6BQ, UK
E-mail: ruth.dobson@qmul.ac.uk

Keywords: multiple sclerosis; genetics; ethnicity

Abbreviations: AUC = area under the curve; EBV = Epstein–Barr virus; EPIC = Expression, Proteomics, Imaging, Clinical; G&H = Genes & Health; GWAS = genome-wide association study(ies); HES = hospital episode statistics; HLA = human leukocyte antigen; HWE = Hardy–Weinberg equilibrium; ICD = International Classification of Diseases; IMSGC = International Multiple Sclerosis Genetics Consortium; LD = linkage disequilibrium; MAF = minor allele frequency; MHC = major histocompatibility complex

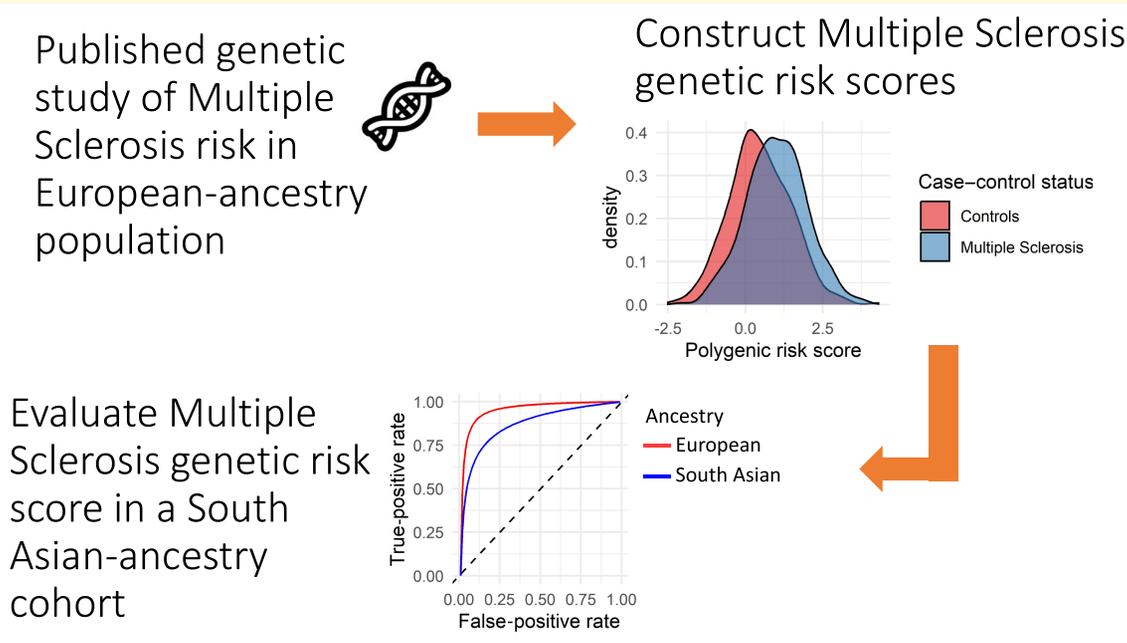
Received July 25, 2022. Revised October 12, 2022. Accepted February 21, 2023. Advance access publication February 22, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of the Guarantors of Brain.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

OR = odds ratio; PCA = principal component analysis; PRS = polygenic risk score(s); SNOMED = Systematized Nomenclature of Medicine; SNP = single nucleotide polymorphism; TOPMed = Trans-Omics for Precision Medicine Program; UCSF = University of California San Francisco; UKB = UK Biobank

Graphical Abstract



Introduction

An individual's risk of developing multiple sclerosis is influenced by common variation across the genome.^{1,2} Multiple sclerosis is a typical complex disease in which the genetic contribution to risk is governed by a large number of susceptibility alleles with individually weak effects. Variation within the major histocompatibility complex (MHC) has the greatest impact on individual risk [odds ratio (OR) associated with DRB1*1501 3.1 and 6.2 for heterozygous and homozygous carriage, respectively].^{2,3} Genome-wide association studies (GWAS) of multiple sclerosis susceptibility have demonstrated at least 200 risk alleles outside the MHC locus, each with a small incremental effect (OR per allele ≤ 1.3).² There is no convincing evidence for monogenic forms of multiple sclerosis in the general population.⁴

Predicting who is likely to develop multiple sclerosis in the future has potential utility for research studies. Accurate disease prediction could facilitate the design of trials for candidate preventive strategies, such as an Epstein-Barr virus (EBV) vaccine or a vitamin D supplementation trial. As multiple sclerosis is a relatively rare disease, such trials will only have the power to demonstrate a risk reduction if the trial population is sufficiently enriched with people at high risk of multiple sclerosis, effectively

increasing the proportion likely to develop the disease.⁵ Furthermore, identifying those at highest risk of disease may allow treatment during the 'prodromal' period, prior to overt clinical manifestations.⁶

Polygenic risk scores (PRS) summarize an individual's cumulative burden of genetic risk alleles to approximate their overall disease risk. Most PRS are calculated by weighting the individual's burden of risk alleles by the estimated effect of each allele on risk—these estimates are usually obtained from GWAS. In two large cohort studies—UK Biobank (UKB) and University of California San Francisco (UCSF) Expression, Proteomics, Imaging, Clinical (EPIC)—PRS have been empirically demonstrated to distinguish multiple sclerosis cases from controls at a population level.⁷⁻⁹

PRS perform poorly in non-European ancestral groups, a phenomenon largely due to differences in linkage disequilibrium (LD) and allele frequencies between populations.¹⁰⁻¹² It is now clear that multiple sclerosis affects individuals of all ethnic backgrounds and that, broadly speaking, the genetic architecture of multiple sclerosis susceptibility overlaps considerably between ancestral groups.¹³⁻²² We therefore sought to evaluate the performance of multiple sclerosis PRS in ~50 000 individuals of South Asian ancestry from the Genes & Health (G&H) cohort to determine the applicability of PRS in this population.

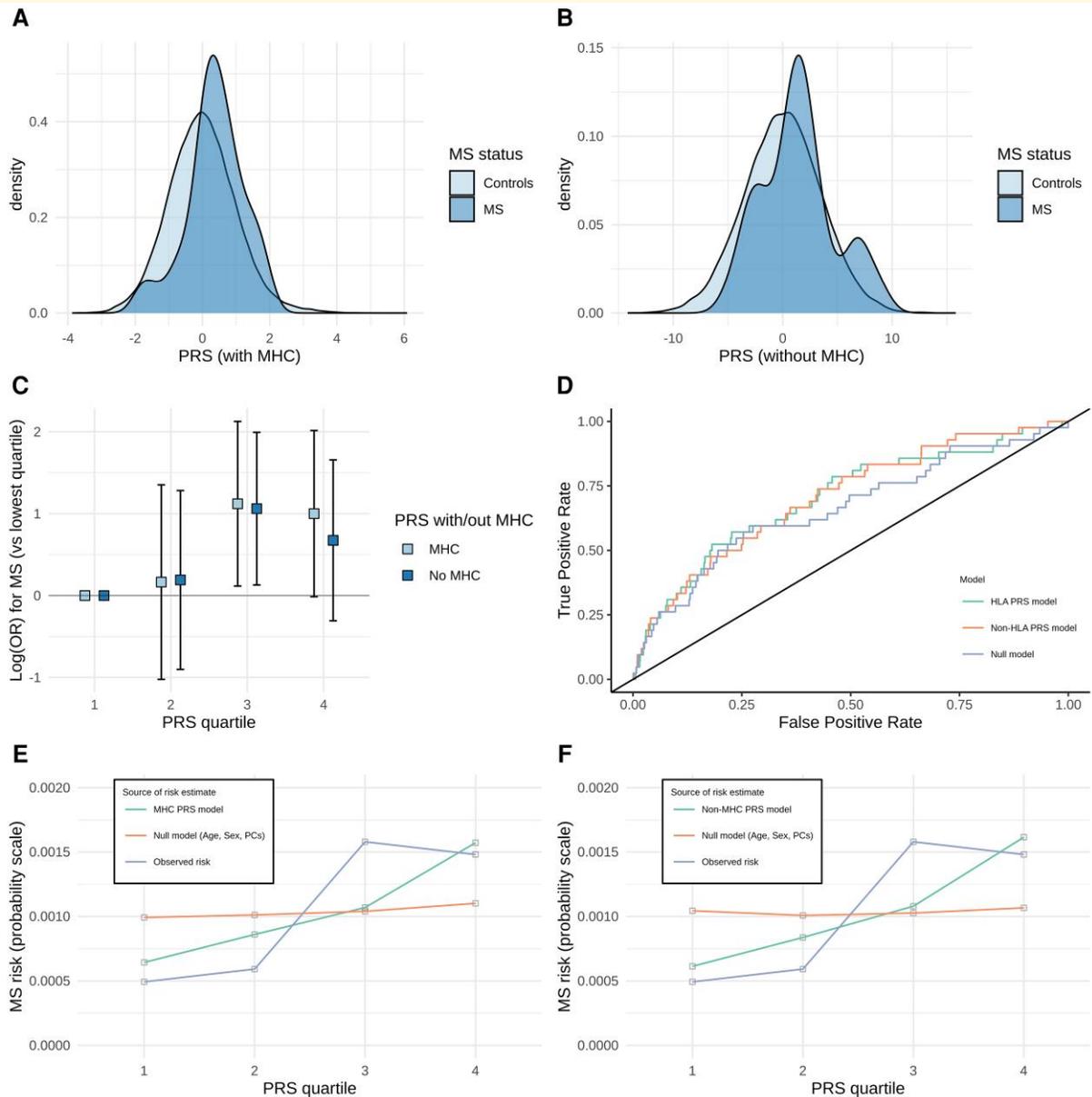


Figure 2 Multiple sclerosis PRS performance in the G&H cohort of South Asian-ancestry individuals. (A and B) Density plots showing the distribution of PRS for PRS with (A) and without (B) the MHC locus in multiple sclerosis cases and controls. (C) Odds ratio quartile plots for individual PRS scores. ORs were calculated relative to the lower quartile. (D) Receiver operating characteristic (ROC) curves for the MHC PRS model, non-MHC model and the null model, with corresponding AUC scores. (E and F) Calibration plots showing the absolute multiple sclerosis disease probabilities (prevalence) for each PRS quartile versus mean fitted probabilities within each quartile from the PRS models. Plots shown for MHC PRS model, non-MHC PRS model, null model and the observed multiple sclerosis risk in each quartile. Odds ratios and AUC values are derived from multivariable logistic regression models

explained $\sim 1.1\%$ of the liability to multiple sclerosis in this cohort (adjusted Nagelkerke's pseudo- R^2 0.011, $P = 0.033$, $N_{\text{SNP}} = 1356$, clumping R^2 0.05, threshold P -value 0.001). The optimal PRS excluding the MHC region ($\text{PRS}_{\text{Non-MHC}}$) performed similarly, explaining $\sim 1.5\%$ of the liability to multiple sclerosis (adjusted Nagelkerke's pseudo- R^2 0.015, $P = 0.015$, $N_{\text{SNP}} = 1965$, clumping R^2 0.4, threshold P -value 0.001). The difference in performance of the $\text{PRS}_{\text{Non-MHC}}$ and PRS_{MHC} was not statistically significant

(likelihood ratio P -value = 1). PRS using variants only lying within the MHC did not correlate with multiple sclerosis disease status ($P = 0.19$).

The predicted risk of multiple sclerosis based on PRS was reasonably well-calibrated to absolute risk (Fig. 2C). Individuals in the top 25% of PRS_{MHC} were nominally more likely to have multiple sclerosis than those in the lowest 25% (OR 2.72, 95% CI 0.99–7.50), although our statistical confidence in this result is tempered by the small number of

poorly in this setting, it does still have some predictive power, consistent with significant overlap in the genetic architecture of multiple sclerosis risk between populations.¹³⁻²²

The lower predictive power of multiple sclerosis PRS we report in an ancestrally South Asian cohort is likely driven by differences in the minor allele frequency of variants and LD structures between European and South Asian populations, rather than due to differences in causal variants.²⁸ If variants included in the PRS are not causal themselves but tag causal variants in Europeans, it does not follow that they will tag the causal variant in other populations, diminishing the accuracy of the score. Previous genetic analyses of multiple sclerosis risk in non-European populations—including small studies of South Asian populations—argue that, broadly speaking, the genetic architecture of multiple sclerosis risk between populations is highly correlated.^{14,16,17,29} Our finding that a European multiple sclerosis PRS has some accuracy in a South Asian cohort, but less so than in Europeans, is entirely consistent with this view.

It is notable that the inclusion of the MHC locus did not improve the PRS in the South Asian cohort. This result could be due to limited statistical power, different causal human leukocyte antigen (HLA) alleles and/or poor tagging of causal HLA alleles by the European GWAS variants. It is important to note that available data suggest that the major HLA risk alleles in Europeans have similar effects in South Asians, and so in our view, it is primarily differences in LD (in addition to the limited case numbers) that drive this unexpected result in the cohort, as well as the statistical imprecision of the effect estimates due to the small number of cases in G&H. Larger studies are required to clarify whether this is merely a power issue.

These results should be interpreted with some degree of caution given the relatively small number of multiple sclerosis cases in the G&H cohort (and the resulting wide confidence intervals), the potential inaccuracies of using electronic health records to ascertain cases (including the possibility of missed cases) and the lack of an external validation cohort. Due to the number of multiple sclerosis cases in G&H, we fitted and evaluated the PRS on the same dataset, which increases the risk of overfitting and therefore may produce an inflated estimate of how well the PRS models disease risk in the population. Furthermore, while we aim to compare PRS performance in UKB and G&H, it is important to note that these cohorts were genotyped on different chips and imputed with different panels (TOPMed versus Haplotype Reference Consortium).^{25,30} Therefore, although we use the same external reference panel to perform LD clumping, the SNPs included in the PRS for any given set of clumping-and-thresholding parameters are not identical between cohorts. The mean age in the G&H cohort is also less than that in UKB, raising the possibility of individuals in the G&H control group going on to develop multiple sclerosis in the future. We aimed to mitigate the effect of sample size by sampling the UKB dataset to an equivalent size.

Given the potential uses of a multiple sclerosis PRS in both clinical care and trial design, the limited cross-ancestry transferability of European-derived PRS is concerning and may reinforce pre-existing health inequalities between different ethnic and ancestral groups. Although advances in statistical methods for applying PRS across populations are likely to enhance transferability,^{11,31} there is an unmet need for ancestrally diverse GWAS of multiple sclerosis risk to ensure that genetics can play a useful role in risk stratification.

Acknowledgements

B.M.J. and R.D. conceived the study. The primary analysis was performed by J.B. and B.M.J. The initial manuscript was drafted by J.B. All authors provided input into critical revisions of the manuscript prior to submission. B.M.J. and J.B. had full access to all the data in the study and take responsibility for the integrity of the data and the accuracy of the data analysis. We thank S.H. for his help with constructing the lay video abstract. We thank Social Action for Health, Centre of The Cell, members of our Community Advisory Group and staff who have recruited and collected data from volunteers. We thank the NIHR National Biosample Centre (UK Biocentre), the Social Genetic and Developmental Psychiatry Centre (King's College London), Wellcome Sanger Institute and Broad Institute for sample processing, genotyping, sequencing and variant annotation. We thank Barts Health NHS Trust, NHS Clinical Commissioning Groups (City and Hackney, Waltham Forest, Tower Hamlets, Newham, Redbridge, Havering, Barking and Dagenham), East London NHS Foundation Trust, Bradford Teaching Hospitals NHS Foundation Trust, Public Health England (especially David Wyllie), Discovery Data Service/Endeavour Health Charitable Trust (especially David Stables) and NHS Digital—for their GDPR-compliant data sharing backed by individual written informed consent. Most of all we thank all of the volunteers participating in G&H.

Funding

This work was performed at the Preventive Neurology Unit, which is funded by the Barts Charity. B.M.J. is supported by a Medical Research Council Clinical Research Training Fellowship (Grant Reference MR/V028766/1) which is co-funded by the UK Multiple Sclerosis Society. Genes & Health is/has recently been core-funded by Wellcome (WT102627, WT210561), the Medical Research Council (UK) (M009017, MR/X009777/1), Higher Education Funding Council for England Catalyst, Barts Charity (845/1796), Health Data Research UK (for London substantive site) and research delivery support from the National Health Service National Institute for Health Research Clinical Research Network (North Thames). Genes & Health is/has recently been funded by Alnylam Pharmaceuticals, Genomics PLC and a Life Sciences

26. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: Rising to the challenge of larger and richer datasets. *Gigascience*. 2015;4:7.
27. King, Butcher & Zalewski. *Apocrita-high performance computing cluster for Queen Mary University of London*. Zenodo; 2017.
28. Wang Y, Guo J, Ni G, Yang J, Visscher PM, Yengo L. Theoretical and empirical quantification of the accuracy of polygenic scores in ancestry divergent populations. *Nat Commun*. 2020;11:3865.
29. Pandit L, Ban M, Sawcer S, *et al*. Evaluation of the established non-MHC multiple sclerosis loci in an Indian population. *Mult Scler*. 2011;17:139-143.
30. Luo Y, Kanai M, Choi W, *et al*. A high-resolution HLA reference panel capturing global population diversity enables multi-ancestry fine-mapping in HIV host response. *Nat Genet*. 2021;53:1504-1516.
31. Ruan Y, Lin Y-F, Feng Y-CA, *et al*. Improving polygenic prediction in ancestrally diverse populations. *Nat Genet*. 2022;54:573-580.