

# 1 The effect of sample size on polygenic hazard models for prostate cancer

2

3 Roshan A. Karunamuni<sup>1</sup>, Minh-Phuong Huynh-Le<sup>1</sup>, Chun C. Fan<sup>2</sup>, Rosalind A.

4 Eeles<sup>3,4</sup>, Douglas F. Easton<sup>5</sup>, ZSofia Kote-Jarai<sup>3</sup>, Ali Amin Al Olama<sup>5,6</sup>, Sara

5 Benlloch Garcia<sup>5</sup>, Kenneth Muir<sup>7,8</sup>, Henrik Gronberg<sup>9</sup>, Fredrik Wiklund<sup>9</sup>, Markus

6 Aly<sup>9,10,11</sup>, Johanna Schleutker<sup>12,13</sup>, Csilla Sipeky<sup>12</sup>, Teuvo LJ Tammela<sup>14,15</sup>, Børge

7 G. Nordestgaard<sup>16,17</sup>, Tim J. Key<sup>18</sup>, Ruth C. Travis<sup>18</sup>, David E. Neal<sup>19,20,21</sup>, Jenny

8 L. Donovan<sup>22</sup>, Freddie C. Hamdy<sup>23,24</sup>, Paul Pharoah<sup>25</sup>, Nora Pashayan<sup>26,25,27</sup>,

9 Kay-Tee Khaw<sup>28</sup>, Stephen N. Thibodeau<sup>29</sup>, Shannon K. McDonnell<sup>30</sup>, Daniel J.

10 Schaid<sup>30</sup>, Christiane Maier<sup>31</sup>, Walther Vogel<sup>32</sup>, Manuel Luedeke<sup>31</sup>, Kathleen

11 Herkommer<sup>33</sup>, Adam S. Kibel<sup>34</sup>, Cezary Cybulski<sup>35</sup>, Dominika Wokolorczyk<sup>35</sup>,

12 Wojciech Kluzniak<sup>35</sup>, Lisa Cannon-Albright<sup>36,37</sup>, Hermann Brenner<sup>38,39,40</sup>, Ben

13 Schöttker<sup>41,42</sup>, Bernd Holleczer<sup>43,44</sup>, Jong Y. Park<sup>45</sup>, Thomas A. Sellers<sup>45</sup>, Hui-Yi

14 Lin<sup>46</sup>, Chavdar Slavov<sup>47</sup>, Radka Kaneva<sup>48</sup>, Vanio Mitev<sup>48</sup>, Jyotsna Batra<sup>49,50</sup>,

15 Judith A. Clements<sup>51,52</sup>, Amanda Spurdle<sup>53</sup>, Australian Prostate Cancer

16 BioResource (APCB)<sup>51</sup>, Manuel R. Teixeira<sup>54,55</sup>, Paula Paulo<sup>54,56</sup>, Sofia Maia<sup>54,56</sup>,

17 Hardev Pandha<sup>57</sup>, Agnieszka Michael<sup>57</sup>, Ian G. Mills<sup>58,59</sup>, Ole A. Andreassen<sup>60</sup>,

18 Anders M. Dale<sup>61,62,63</sup>, Tyler M. Seibert<sup>1</sup>, The PRACTICAL Consortium<sup>^</sup>

19 <sup>1</sup>Department of Radiation Medicine and Applied Sciences, University of California  
20 San Diego, La Jolla, CA, USA

21 <sup>2</sup>Healthytix, 4747 Executive Dr. Suite 820, San Diego, CA, USA

22 <sup>3</sup>The Institute of Cancer Research, London, SM2 5NG, UK

23 <sup>4</sup>Royal Marsden NHS Foundation Trust, London, SW3 6JJ, UK

24 <sup>5</sup>Centre for Cancer Genetic Epidemiology, Department of Public Health and  
25 Primary Care, University of Cambridge, Strangeways Research Laboratory,  
26 Cambridge CB1 8RN, UK

27 <sup>6</sup>University of Cambridge, Department of Clinical Neurosciences, Stroke  
28 Research Group, R3, Box 83, Cambridge Biomedical Campus, Cambridge CB2

- 29 0QQ, UK
- 30 <sup>7</sup>Division of Population Health, Health Services Research and Primary Care,  
31 University of Manchester, Oxford Road, Manchester, M13 9PL, UK
- 32 <sup>8</sup>Warwick Medical School, University of Warwick, Coventry, UK
- 33 <sup>9</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institute, SE-  
34 171 77 Stockholm, Sweden
- 35 <sup>10</sup>Department of Molecular Medicine and Surgery, Karolinska Institute, SE-171  
36 77 Stockholm, Sweden
- 37 <sup>11</sup>Department of Urology, Karolinska University Hospital, Stockholm, Sweden
- 38 <sup>12</sup>Institute of Biomedicine, Kiinamyllynkatu 10, FI-20014 University of Turku,  
39 Finland
- 40 <sup>13</sup>Department of Medical Genetics, Genomics, Laboratory Division, Turku  
41 University Hospital, PO Box 52, 20521 Turku, Finland
- 42 <sup>14</sup>Faculty of Medicine and Health Technology, Prostate Cancer Research Center,  
43 FI-33014 Tampere University, Finland
- 44 <sup>15</sup>Department of Urology, Tampere University Hospital, Tampere, Finland
- 45 <sup>16</sup>Faculty of Health and Medical Sciences, University of Copenhagen, 2200  
46 Copenhagen, Denmark
- 47 <sup>17</sup>Department of Clinical Biochemistry, Herlev and Gentofte Hospital,  
48 Copenhagen University Hospital, Herlev, 2200 Copenhagen, Denmark
- 49 <sup>18</sup>Cancer Epidemiology Unit, Nuffield Department of Population Health, University  
50 of Oxford, Oxford, OX3 7LF, UK
- 51 <sup>19</sup>Nuffield Department of Surgical Sciences, University of Oxford, Room 6603,  
52 Level 6, John Radcliffe Hospital, Headley Way, Headington, Oxford, OX3 9DU,  
53 UK
- 54 <sup>20</sup>University of Cambridge, Department of Oncology, Box 279, Addenbrooke's  
55 Hospital, Hills Road, Cambridge CB2 0QQ, UK
- 56 <sup>21</sup>Cancer Research UK, Cambridge Research Institute, Li Ka Shing Centre,  
57 Cambridge UK
- 58 <sup>22</sup>School of Social and Community Medicine, University of Bristol, Canynge Hall,  
59 39 Whatley Road, Bristol, BS8 2PS, UK
- 60 <sup>23</sup>Nuffield Department of Surgical Sciences, University of Oxford, Oxford, OX1  
61 2JD, UK
- 62 <sup>24</sup>Faculty of Medical Science, University of Oxford, John Radcliffe Hospital,  
63 Oxford, UK
- 64 <sup>25</sup>Centre for Cancer Genetic Epidemiology, Department of Oncology, University  
65 of Cambridge, Strangeways Laboratory, Worts Causeway, Cambridge, CB1  
66 8RN, UK
- 67 <sup>26</sup>University College London, Department of Applied Health Research,  
68 London, UK
- 69 <sup>27</sup>Department of Applied Health Research, University College London, London,  
70 WC1E 7HB, UK
- 71 <sup>28</sup>Clinical Gerontology Unit, University of Cambridge, Cambridge, CB2 2QQ, UK
- 72 <sup>29</sup>Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester,  
73 MN 55905, USA
- 74 <sup>30</sup>Division of Biomedical Statistics & Informatics, Mayo Clinic, Rochester, MN

- 75 55905, USA
- 76 <sup>31</sup>Humangenetik Tuebingen, Paul-Ehrlich-Str 23, D-72076 Tuebingen
- 77 <sup>32</sup>Institute for Human Genetics, University Hospital Ulm, 89075 Ulm, Germany
- 78 <sup>33</sup>Technical University of Munich, School of Medicine, Klinikum rechts der Isar,
- 79 Department of Urology
- 80 <sup>34</sup>Division of Urologic Surgery, Brigham and Womens Hospital, 75 Francis Street,
- 81 Boston, MA 02115, USA
- 82 <sup>35</sup>International Hereditary Cancer Center, Department of Genetics and Pathology,
- 83 Pomeranian Medical University, 70-115 Szczecin, Poland
- 84 <sup>36</sup>Division of Genetic Epidemiology, Department of Medicine, University of Utah
- 85 School of Medicine, Salt Lake City, Utah 84112, USA
- 86 <sup>37</sup>George E. Wahlen Department of Veterans Affairs Medical Center, Salt Lake
- 87 City, Utah 84148, USA
- 88 <sup>38</sup>Division of Clinical Epidemiology and Aging Research, German Cancer
- 89 Research Center (DKFZ), D-69120, Heidelberg, Germany
- 90 <sup>39</sup>German Cancer Consortium (DKTK), German Cancer Research Center
- 91 (DKFZ), D-69120 Heidelberg, Germany
- 92 <sup>40</sup>Division of Preventive Oncology, German Cancer Research Center (DKFZ) and
- 93 National Center for Tumor Diseases (NCT), Im Neuenheimer Feld 460
- 94 69120 Heidelberg, Germany
- 95 <sup>41</sup>Division of Clinical Epidemiology and Aging Research, German Cancer
- 96 Research Center (DKFZ), D-69120 Heidelberg, Germany
- 97 <sup>42</sup>Network Aging Research, University of Heidelberg, Heidelberg, Germany
- 98 <sup>43</sup>Saarland Cancer Registry, D-66119 Saarbrücken, Germany
- 99 <sup>44</sup>Division of Clinical Epidemiology and Aging Research, German Cancer
- 100 Research Center (DKFZ), Heidelberg, Germany
- 101 <sup>45</sup>Department of Cancer Epidemiology, Moffitt Cancer Center, 12902 Magnolia
- 102 Drive, Tampa, FL 33612, USA
- 103 <sup>46</sup>School of Public Health, Louisiana State University Health Sciences Center,
- 104 New Orleans, LA 70112, USA
- 105 <sup>47</sup>Department of Urology and Alexandrovska University Hospital, Medical
- 106 University of Sofia, 1431 Sofia, Bulgaria
- 107 <sup>48</sup>Molecular Medicine Center, Department of Medical Chemistry and
- 108 Biochemistry, Medical University of Sofia, Sofia, 2 Zdrave Str., 1431 Sofia,
- 109 Bulgaria
- 110 <sup>49</sup>Institute of Health and Biomedical Innovation and School of Biomedical
- 111 Sciences, Queensland University of Technology, Brisbane, QLD 4059, Australia
- 112 <sup>50</sup>Australian Prostate Cancer Research Centre-Qld, Translational Research
- 113 Institute, Brisbane, Queensland 4102, Australia
- 114 <sup>51</sup>Australian Prostate Cancer Research Centre-Qld, Institute of Health and
- 115 Biomedical Innovation and School of Biomedical Science, Queensland University
- 116 of Technology, Brisbane QLD 4059, Australia
- 117 <sup>52</sup>Translational Research Institute, Brisbane, Queensland 4102, Australia
- 118 <sup>53</sup>Molecular Cancer Epidemiology Laboratory, QIMR Berghofer Institute of
- 119 Medical Research, Brisbane, Australia
- 120 <sup>54</sup>Department of Genetics, Portuguese Oncology Institute of Porto (IPO-Porto),

4200-072 Porto, Portugal

<sup>55</sup>Biomedical Sciences Institute (ICBAS), University of Porto, 4050-313 Porto, Portugal

<sup>56</sup>Cancer Genetics Group, IPO-Porto Research Center (CI-IPOP), Portuguese Oncology Institute of Porto (IPO-Porto), Porto, Portugal

<sup>57</sup>The University of Surrey, Guildford, Surrey, GU2 7XH, UK

<sup>58</sup>Center for Cancer Research and Cell Biology, Queen's University of Belfast, Belfast, UK

<sup>59</sup>Nuffield Department of Surgical Sciences, John Radcliffe Hospital, University of Oxford, Oxford, UK

<sup>60</sup>NORMENT, KG Jebsen Centre, Oslo University Hospital and University of Oslo, Oslo, Norway

<sup>61</sup>Department of Radiology, University of California San Diego, La Jolla, CA, USA

<sup>62</sup>Department of Cognitive Science, University of California San Diego, La Jolla, CA, USA

<sup>63</sup>Department of Neurosciences, University of California San Diego, La Jolla, CA, USA

\*Corresponding Author:

E-mail: [rakarunamuni@ucsd.edu](mailto:rakarunamuni@ucsd.edu)(RK), [tseibert@ucsd.edu](mailto:tseibert@ucsd.edu)(TM)

<sup>^</sup> Membership of The PRACTIAL Consortium is provided in the Supporting Information.

# **Abstract**

We aimed to determine the effect of sample size on performance of polygenic hazard score (PHS) models in predicting the age at onset of prostate cancer. Age and genotypes were obtained for 40,861 men from the PRACTICAL consortium. The dataset included 201,590 SNPs per subject, and was split into training (34,444 samples) and testing (6,417 samples) sets. Two PHS model-building strategies were investigated. Established-SNP model considered 65 SNPs that had been associated with prostate cancer in the literature. A stepwise SNP selection was used to develop Discovery-SNP models. The performance of each PHS model was calculated for random sizes of the training set (1 to 30 thousand). The performance of a representative Established-SNP model was estimated for random sizes of the testing set (0.5 to 6 thousand). Mean  $HR_{98/50}$  (hazard ratio of top 2% to the average in the test set) of the Established-SNP model increased from 1.73[95%CI: 1.69-1.77] to 2.41[2.40-2.43] when the number of training samples was increased from 1 to 30 thousand. The corresponding  $HR_{98/50}$  of the Discovery-SNP model increased from 1.05[0.93-1.18] to 2.19[2.16-2.23].  $HR_{98/50}$  of a representative Established-SNP model using testing set sample sizes of 0.6 and 6 thousand observations were 1.78[1.70-1.85] and 1.73[1.71-1.76], respectively. We estimate that a study population of 20 to 30 thousand men is required to develop Discovery-SNP PHS models for prostate cancer. The required sample size could be reduced to 10 thousand samples, if a set of SNPs associated with the disease has already been established.

## Author summary

Polygenic hazard scores represent a recent advancement in polygenic prediction to model the age of onset of various diseases, such as Alzheimer's disease or prostate cancer. These scores accumulate small effect sizes from several tens of genetic variants and can be used to establish an individual's risk of experiencing an event relative to a control population across time. The largest barrier to the development of polygenic hazard scores is the large number of study subjects needed to develop the underlying models. We sought to understand the effect of varying the total number of samples on the performance of a polygenic hazard score in the context of prostate cancer. We found that the performance of the score did not appreciably change beyond 20 to 30 thousand observations when developing the model from scratch. However, when the discovery of the genetic variants can be borrowed from those already identified in the literature to be associated with the disease, the required number of samples is reduced to 10 thousand with no appreciable detriment in performance. We hope that these results can guide the design of future studies of polygenic scores in other diseases and demonstrate the importance of genome-wide association studies.

## Introduction

Polygenic prediction models have been studied extensively for several diseases such as prostate cancer[1], breast cancer[2], type 2 diabetes[3], dementia[4], and atherosclerosis[5]. Polygenic scores in the context of survival models are a more recent advancement in the field, but have been garnering interest in the prediction of age at onset of Alzheimer's disease[6] and prostate cancer[7]. The steady increase in genetic testing[8,9], both in public and clinical domains, suggests that survival models could be applied to new diseases. The largest obstacle to the development of these models is the large number of study subjects, often in the tens of thousands[8], which are required for robust training and testing.

Our aim was to quantify the effect of sample size on the performance of a polygenic survival model. This was explored through a specific disease condition that is expected to be representative, namely the prediction of age of onset in prostate cancer. We investigated two potential model development strategies. For the 'Established-SNP' model, we selected single-nucleotide polymorphisms (SNPs) that had previously been shown to be associated with prostate cancer, and simply estimated the coefficients for these SNPs in a Cox proportional hazards framework. For the 'Discovery-SNP' model, we implemented the SNP selection technique described by Seibert *et al.*[7] to identify SNPs in our genotyping data for inclusion in the Cox proportional hazards framework. The Established-SNP and Discovery-SNP represent two strategies that researchers could employ to build a polygenic survival model. In order to simulate samples of

different sizes, we randomly sampled our training and testing sets. The results of this work will help inform the design of future studies to develop polygenic survival models for other diseases.

## Results

### Established- vs. Discovery-SNP model performance

Histogram comparisons of performance metrics of Established (EST) and Discovery (DIS) SNP models are illustrated in Figure 1. The performance metrics are shown for 50 random samplings of the training set using a sample size of 30 thousand total observations. Qualitatively, there appears to be more variability in performance metrics associated with the Discovery process.

### Coefficients of Established-SNP model

The mean coefficients for the 65 SNPs used in the Established-SNP model are plotted in Figure 2.

### Effect of training set sample size on performance

Box plots of the performance metrics of the Established-SNP and Discovery-SNP models for random samples of the training set are shown in Figure 3 and Figure 4, respectively. The mean values of  $HR_{98/50}$ ,  $HR_{20/50}$ ,  $HR_{98/20}$ ,  $HR_{80/20}$ , z-score, and beta using a random training sample of 1 thousand total observations in the Established-SNP model were 1.73 [95% CI: 1.69-1.76], 0.71 [0.71-0.73], 2.42 [2.35-2.50], 1.96 [1.92-2.01], 9.92 [9.57-10.28], and 0.45 [0.43-



0.47] respectively. The corresponding values using a random training sample of 30 thousand total observations were 2.41 [95% CI: 2.40-2.43], 0.60 [0.60-0.60], 4.04 [4.02-4.07], 2.86 [2.84-2.87], 15.1 [15.04-15.16], and 1.18 [1.17-1.18] respectively.

The mean values of  $HR_{98/50}$ ,  $HR_{20/50}$ ,  $HR_{98/20}$ ,  $HR_{80/20}$ , z-score, and beta using a random training sample of 1 thousand total observations in the Discovery-SNP model were 1.05 [0.93-1.18], 0.98 [0.89-1.07], 1.07 [0.91-1.24], 1.08 [0.91-1.24], 1.06 [-1.20-3.31], and 0.17 [-0.23-0.65] respectively. The corresponding performance values using a training sample size of 30 thousand observations were 2.20 [2.16-2.23], 1.60 [1.59-1.62], 3.47 [3.39-3.56], 2.53 [2.49-2.58], 13.19 [12.96-13.41], and 0.87 [0.85-0.89] respectively.

#### Effect of testing set sample size on performance

Box plots of the performance metrics of the representative Established-SNP model for random samples of the testing set are shown in Figure 5. The mean values of  $HR_{98/50}$ ,  $HR_{20/50}$ ,  $HR_{98/20}$ ,  $HR_{80/20}$ , z-score, and beta using a random testing sample of 0.5 thousand total observations in the representative Established-SNP model were 1.78 [1.71-1.85], 0.73 [0.71-0.74], 2.50 [2.33-2.66], 1.99 [1.89-2.09], 3.82 [3.57-4.08], and 0.76 [0.70-0.82] respectively. The corresponding values using a testing sample of 6 thousand observations were: 1.73 [1.72-1.76], 0.73 [0.72-0.73], 2.39 [2.34-2.44], 1.93 [1.90-1.96], 13.07 [12.80-13.32], and 0.74 [0.72-0.76] respectively.

## Discussion

We identified several trends in the effect of training and testing sample size on the performance of PHS models in predicting the age of onset of prostate cancer using SNP genetic variants. When using SNPs that had already been associated with prostate cancer risk, our analysis suggests that very little improvement in performance can be achieved once the training sets becomes larger than 10 to 15 thousand observations. When attempting to discover SNPs, a similar plateau in performance was observed from training sets larger than 20 to 25 thousand observations. Apart from z-scores, the performance metrics of the chosen Cox proportional hazards model did not vary with testing sample size. However, we did observe that the distribution of performance metrics narrows until a testing sample size of 3 to 4 thousand observations, after which the distribution remains relatively stable.

Our results may be used to inform researchers on the approximate number of subjects needed to develop PHS models to predict the age of onset of diseases using SNP counts. A dataset of 20 thousand observations may be the minimum needed to accurately estimate the PHS coefficients of SNPs that have been previously discovered in the setting of a logistic model. Such a dataset would allow for the accurate estimation of SNP coefficients as well as the testing of model performance in an independent holdout set. Based on our results, this number would have to be increased to roughly 30 thousand observations if the researchers intend on discovering the SNPs from scratch using the approach described here.

The PHS model developed by Desikan *et al.*[6] to predict age-associated risk of Alzheimer's disease used a training set with roughly 55,000 individuals. A similarly structured model developed by Seibert *et al.*[7] to guide screening for aggressive prostate cancer was developed with roughly 31,000 men. Studies such as these require large investments in time, money, and resources in order to acquire the genetic data needed for the analysis. The results of our analysis help elucidate that the minimum sample size needed to translate this technology to other diseases and processes may be lower than what has been used so far in previous studies. This seems to be particularly true if the researchers use SNPs that have already been discovered and validated as associated with the process of interest.

The results of this study must be considered in the context of its limitations. The list of Established-SNPs was previously selected from a larger dataset that included the sample patients used in the test set in the present study. As such, there is leakage of information from the test set to the development of the Established-SNP model. Therefore, the performance metrics of the Established-SNP model should not be directly compared to those of the Discovery-SNP model, as the values of the former may be inflated.

In addition, we have chosen to focus on only two of countless possible model development schemes. The role of sample size in other development strategies—such as regularized Cox proportional models, parametric survival functions, or random survival forests—is yet to be explored. Finally, the analysis is limited to prostate cancer and to the SNPs on the iCOGS array. Future work

will include SNPs imputed from 1000 Genomes[13]. Such an analysis was not performed for this first study to limit computation time for bootstrap analyses and to avoid uncertainty due to imputation.

In conclusion, we have studied the effect of sample size on the performance of PHS models to study the association between SNPs and the age at onset of prostate cancer. We have determined that models require roughly 20 to 30 thousand samples before their performance would not be improved greatly by expansion of the training set. Using SNPs that have already been established in the literature may help reduce the number of training samples required to reach this performance plateau by almost 10 thousand samples.

## **Materials and Methods**

### Training and testing set

As previously described[7], we obtained genotype and age data from 21 studies included in the Prostate Cancer Association Group to Investigate Cancer Associated Alterations in the Genome (PRACTICAL) consortium. We analyzed data from 40,861 men consisting of 20,551 individuals with prostate cancer and 20,310 individuals without. For analysis, the age for each man was recorded as either their age at prostate cancer diagnosis (cases) or at interview (controls). Genotype data for 201,590 SNPs were also available for analysis. The genotype data had been assayed using a custom iCOGS chip (Illumina, San Diego, CA) the details for which are elaborated elsewhere[10]. The sample was split into training (34,444 men) and testing (6,417 men) sets. The testing set was selected

using men who were enrolled in the Prostate testing for cancer and Treatment ( ProtecT[11]) trial. ProtecT (ClinicalTrials.gov: NCT02044172) is a large, multicenter trial within the United Kingdom which aims to investigate the effectiveness of treatments for localized prostate cancer. The ProtecT study group was chosen for testing as it represented a well-characterized group of individuals that had been used for measuring testing performance for our earlier work. The Data Availability Statement describing how readers can gain access to the PRACTICAL dataset is provided in the Supplementary Information.

### Established-SNP model

A list of 65 SNPs[12] was chosen to represent those on the iCOGS array that had been published as associated with prostate cancer. The coefficients of the SNPs within the Established-SNP model were then estimated using the “coxphfit” function in MATLAB (Mathworks, Natwick, MA). Prior to parameter estimation, missing SNP data were replaced by mean imputation. It should be noted that the 65 SNPs used were discovered, in large part, using the data presently defined as the test set. The effect allele for all 65 SNPs was defined as “A” to simplify analysis.

### Discovery-SNP model

SNPs with call rates less than 95% were removed from the selection process. For every SNP, a trend test was used to check for associations between SNP count and the binary classification of individuals with or without prostate

cancer. The SNP selection pool was then reduced to those whose trend test p-value was less  $1 \times 10^{-6}$ . In order of increasing p-value, each SNP was tested in a multiple logistic regression model for association with the binary classification of men as with or without prostate cancer, after adjusting for age, six principal components based upon genetic ancestry, and previously selected SNPs. If the p-value of the coefficient of the tested SNP was less than  $1 \times 10^{-6}$ , it was selected for the final Cox proportional hazard model estimation. The coefficients of the selected SNP pool within the Discovery-SNP model were estimated as previously described[7].

### Polygenic Hazard Score (PHS)

The polygenic hazard score (PHS) for each of the Established-SNP and Discovery-SNP models was calculated as the linear product of the coefficients of the SNPs used in the model and the corresponding patient genotype counts[6,7].

### PHS performance metrics

Several performance metrics for PHS models were investigated, and are described in Table 1. In each case, the PHS for each test subject was calculated as the dot product of SNP coefficients, either Established or Discovery, and SNP counts. A Cox proportional hazards model was then fit using PHS as the sole predictor of age in the test set. The z-score and beta of this Cox proportional hazards model relate to how well PHS was associated with age within the test set. The hazard ratios were calculated as the exponential of the differences in

predicted log-relative hazards of different groups within the test set. The groups were defined using centile cut-points for those controls within the training set whose age was less than 70 years. This list of performance metrics expands on those (z-score and  $HR_{98/50}$ ) that were used in our earlier work[7].

**Table 1.** Performance metrics used in the evaluation of polygenic hazard scores.

| Performance metric | Description  |
|--------------------|--|
| $HR_{98/50}$       | Hazard ratio of the top 2% to the average (30 – 70%) in the test set   |
| $HR_{20/50}$       | Hazard ratio of the bottom 20% to the average (30 – 70%) in the test set                                     |
| $HR_{98/20}$       | Hazard ratio of the top 2% to the bottom 20% in the test set   |
| $HR_{80/20}$       | Hazard ratio of the top 20% to the bottom 20% in the test set.   |
| z-score            | z-score of Cox proportional hazards model using PHS as a sole predictor of age in the test set               |
| beta               | coefficient of PHS in a Cox proportional hazards model using PHS as a sole predictor of age in the test set. |

Random sampling of training set

Random sampling of the training set was performed with replacement while ensuring equal proportions of men with and without prostate cancer. The training set was randomly sampled to include 1, 5, 10, 15, 20, 25, and 30 thousand total observations. Performance of the Established and Discovery-SNP models using random samples of the training data was measured in the entire test set.

#### Random sampling of the testing set

Random sampling of the testing set was performed with replacement while ensuring equal proportion of men with and without prostate cancer. The testing set was randomly sampled to include 0.5, 1, 2, 3, 4, 5 and 6 thousand total observations. Performance in the randomly sampled testing sets was performed using a representative Established-SNP model. The representative model was chosen as that whose parameters were estimated using a training sample size of 30 thousand total observations, and whose performance metrics were the shortest Euclidean distance to the average performance across all Established-SNP models using a training sample size of 30 thousand.

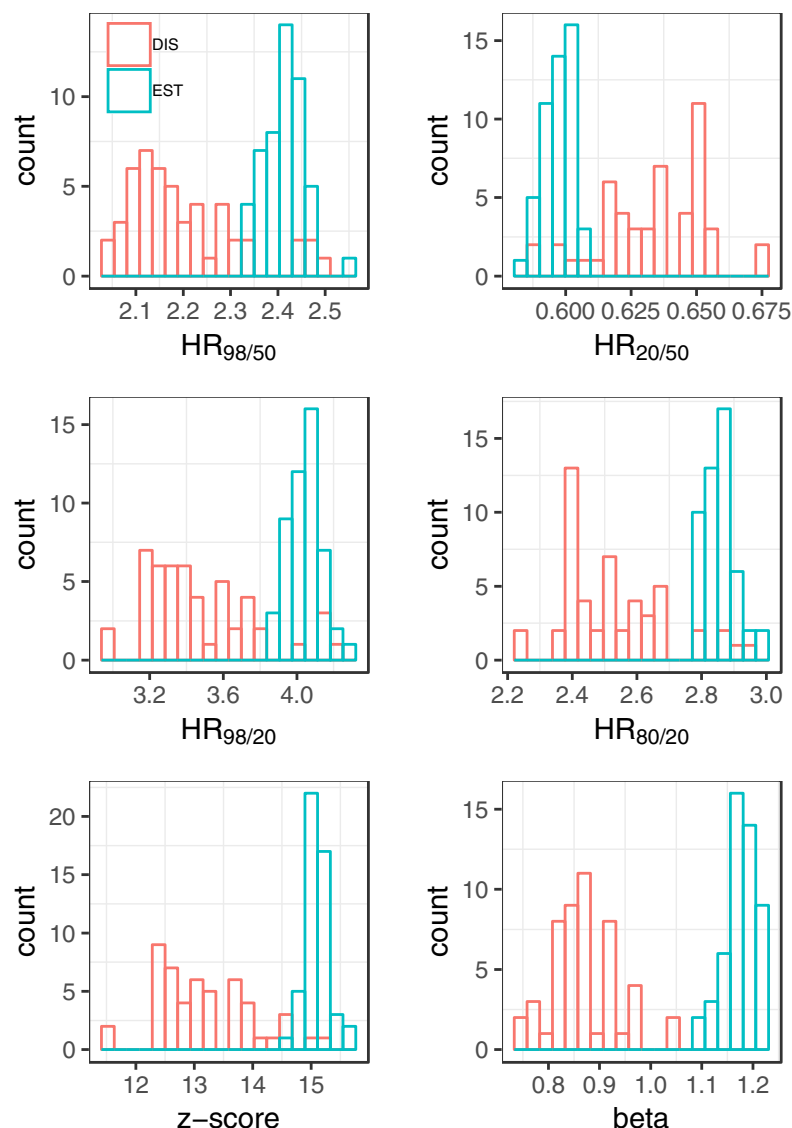


## References

1. Aly M, Wiklund F, Xu J, Isaacs WB, Eklund M, D'Amato M, et al. Polygenic risk score improves prostate cancer risk prediction: Results from the Stockholm-1 cohort study. *Eur Urol*. 2011;60: 21–28. doi:10.1016/j.eururo.2011.01.017
2. Machiela MJ, Chen C, Chanock SJ, Hunter DJ, Kraft P. Evaluation of polygenic risk scores for predicting breast and prostate cancer risk. *Genet Epidemiol*. 2011;514: n/a-n/a. doi:10.1002/gepi.20600
3. Vassy JL, Hivert MF, Porneala B, Dauriz M, Florez JC, Dupuis J, et al. Polygenic type 2 diabetes prediction at the limit of common variant detection. *Diabetes*. 2014;63: 2172–2182. doi:10.2337/db13-1663
4. Marden JR, Walter S, Tchetgen Tchetgen EJ, Kawachi I, Glymour MM. Validation of a polygenic risk score for dementia in black and white individuals. *Brain Behav*. 2014;4: 687–697. doi:10.1002/brb3.248
5. Natarajan P, Young R, Stitzel NO, Padmanabhan S, Baber U, Mehran R, et al. Polygenic risk score identifies subgroup with higher burden of atherosclerosis and greater relative benefit from statin therapy in the primary prevention setting. *Circulation*. 2017;135: 2091–2101. doi:10.1161/CIRCULATIONAHA.116.024436
6. Desikan RS, Fan CC, Wang Y, Schork AJ, Cabral HJ, Cupples LA, et al. Genetic assessment of age-associated Alzheimer disease risk: Development and validation of a polygenic hazard score. *PLoS Med*. 2017;14: 1–17. doi:10.1371/journal.pmed.1002258

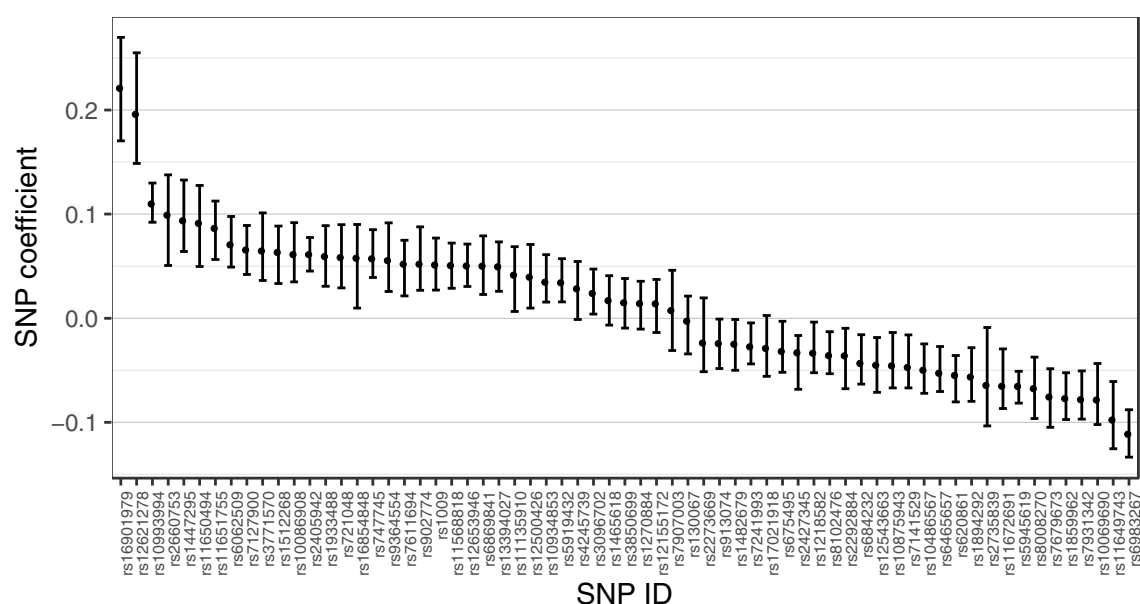
- 420 7. Seibert TM, Fan CC, Wang Y, Zuber V, Karunamuni R, Parsons JK, et al.  
421 Polygenic hazard score to guide screening for aggressive prostate cancer:  
422 Development and validation in large scale cohorts. *BMJ*. 2018;360: 1–7.  
423 doi:10.1136/bmj.j5757
- 424 8. Chatterjee N, Shi J, García-Closas M. Developing and evaluating polygenic  
425 risk prediction models for stratified disease prevention. *Nat Rev Genet*.  
426 Nature Publishing Group; 2016;17: 392–406. doi:10.1038/nrg.2016.27
- 427 9. Torkamani A, Wineinger NE, Topol EJ. The personal and clinical utility of  
428 polygenic risk scores. *Nat Rev Genet*. Springer US; 2018;19: 581–590.  
429 doi:10.1038/s41576-018-0018-x
- 430 10. Eeles RA, Olama AA Al, Benlloch S, Saunders EJ, Leongamornlert DA,  
431 Tymrakiewicz M, et al. Identification of 23 new prostate cancer  
432 susceptibility loci using the iCOGS custom genotyping array. *Nat Genet*.  
433 2013;45: 385–391. doi:10.1038/ng.2560
- 434 11. Lane JA, Donovan JL, Davis M, Walsh E, Dedman D, Down L, et al. Active  
435 monitoring, radical prostatectomy, or radiotherapy for localised prostate  
436 cancer: Study design and diagnostic and baseline results of the ProtecT  
437 randomised phase 3 trial. *Lancet Oncol*. Lane et al. Open Access article  
438 distributed under the terms of CC BY; 2014;15: 1109–1118.  
439 doi:10.1016/S1470-2045(14)70361-4
- 440 12. Szulkin R, Whittington T, Eklund M, Aly M, Eeles RA, Easton D, et al.  
441 Prediction of individual genetic risk to prostate cancer using a polygenic  
442 score. *Prostate*. 2015;75: 1467–1474. doi:10.1002/pros.23037

443 13. Altshuler DL, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A, Clark  
 444 AG, et al. A map of human genome variation from population-scale  
 445 sequencing. *Nature*. 2010;467: 1061–1073. doi:10.1038/nature09534  
 446  
 447



**Figure 1.** Comparison of performance metrics between Established (EST) and Discovery (DIS) SNP models using 50 random samples of the training set using a sample size of 30 thousand. There is more variability with the Discovery process. Established SNPs, though, were discovered using the data in the training set; this circularity is not accounted for in the present study, which focuses on sample size effects.

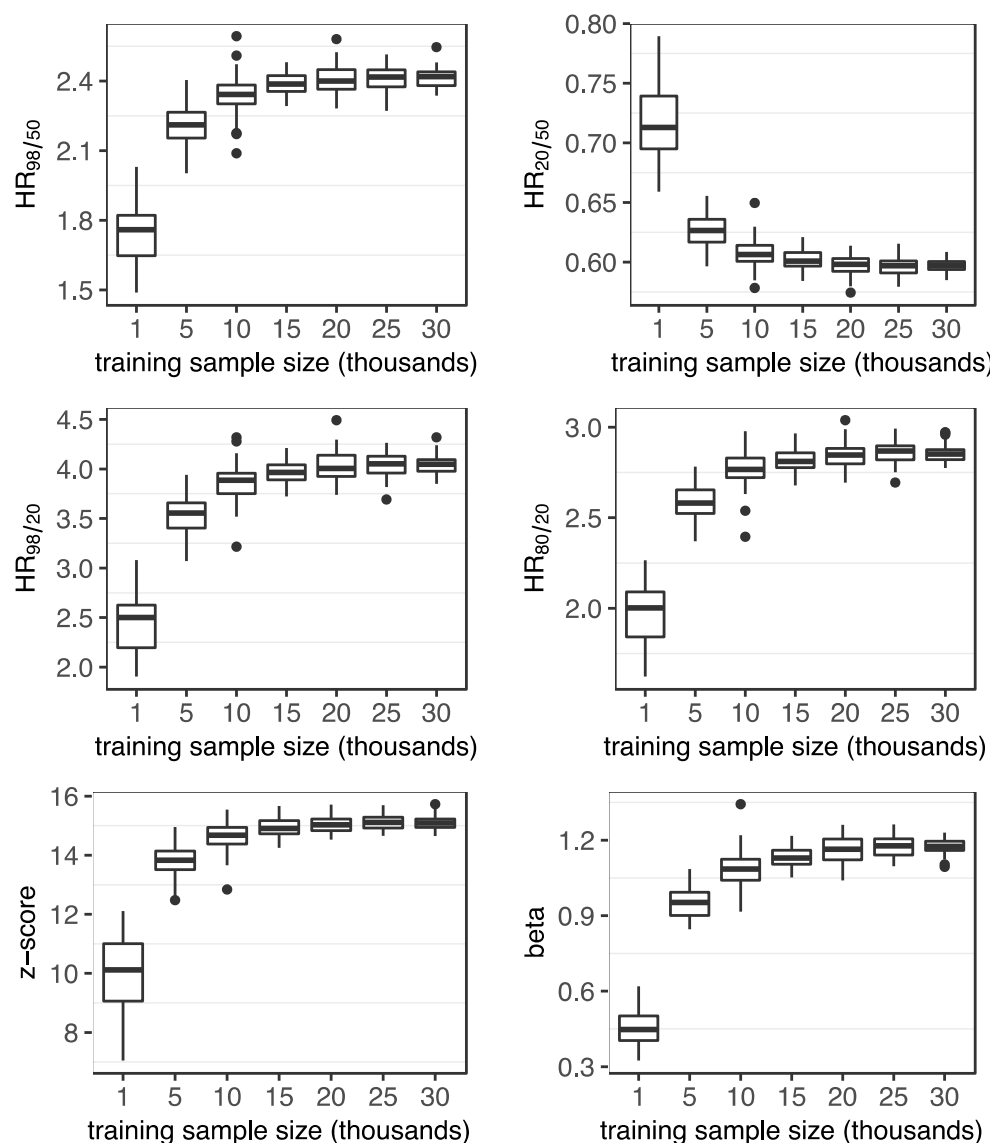
455



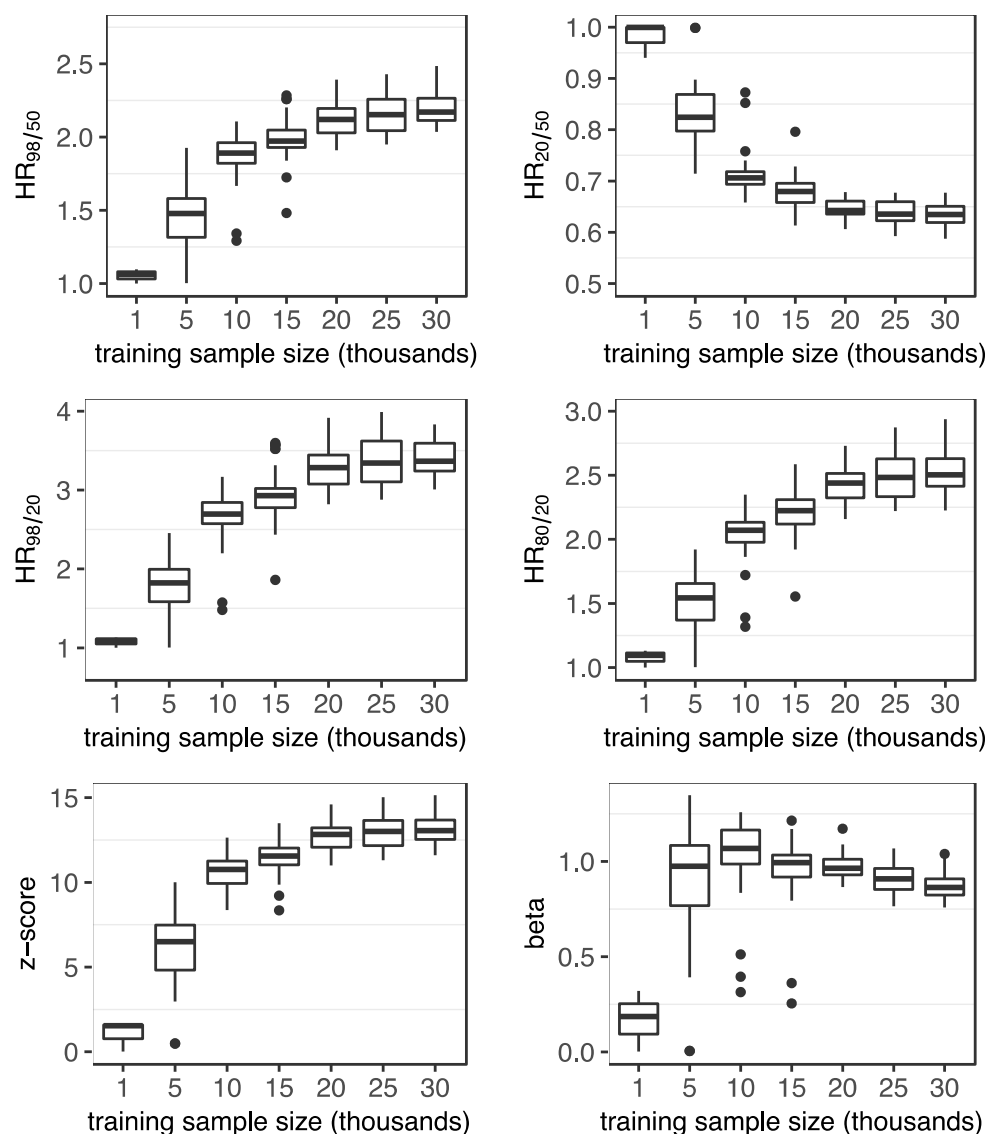
456

457 **Figure 2.** Coefficients of 65 SNPs used in the Established SNP model. Data points represent  
 458 mean values across 50 iterations of a random sample of the training set using a sample size of  
 459 30 thousand total observations. Error bars represent 95% confidence intervals.

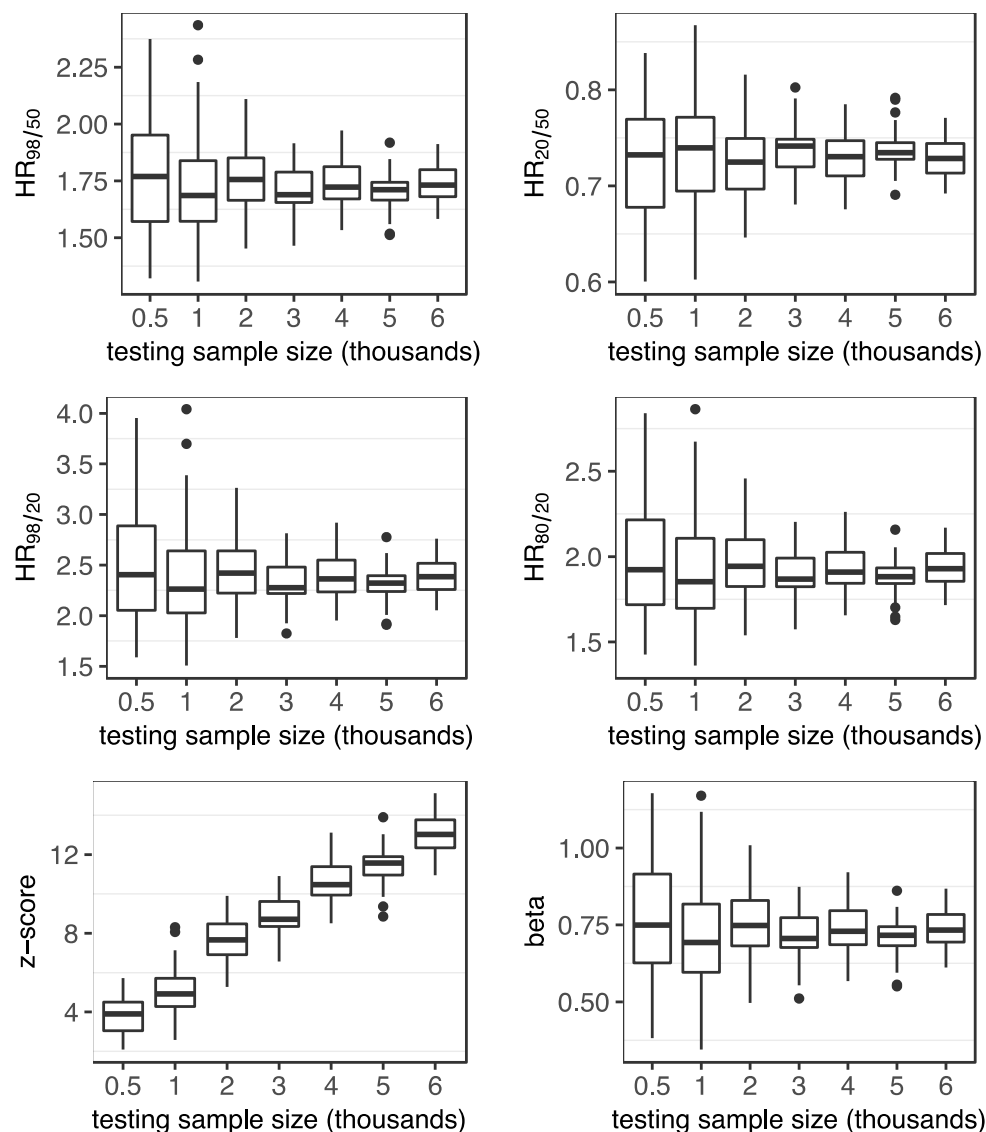
460



**Figure 3.** Performance metrics of Established SNP model. Box plots of performance metrics are shown for random samples of the training set using sample sizes of 1, 5, 10, 15, 20, 25, and 30 thousand total observations. Within each box plot, the horizontal line represents the median and the box extends from the 25<sup>th</sup> to 75<sup>th</sup> percentile.



**Figure 4.** Performance metrics of the Discovery SNP model. Box plots of performance metrics are shown for random samples of the training set using sample sizes of 1, 5, 10, 15, 20, 25, and 30 thousand total observations. Within each box plot, the horizontal line represents the median and the box extends from the 25<sup>th</sup> to 75<sup>th</sup> percentile.



**Figure 5.** Performance as a function of testing sample size. Box plots of performance metrics of the representative Established SNP model in random samples of the testing set from 0.5 to 6 thousand total observations.



478 **Supporting Information Legends**

479 Supporting Information 1. Data Availability Statement details how readers can  
480 obtain the data from the PRACTICAL (Prostate Cancer Association Group to  
481 Investigate Cancer Associated Alterations in the Genome) consortium. The  
482 document also contains the additional authorship, affiliation, and funding sources  
483 for the PRACTICAL consortium.