



Predicting key GWAS outcomes under a point-normal polygenic model

Tian Wu¹, Zipeng Liu², Timothy Shin Heng Mak², Yan Dora Zhang^{2,3}, Pak Chung Sham^{1,2,4}

1. Department of Psychiatry, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China 2. Centre for PanorOmic Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China 3. Department of Statistics & Actuarial Science, Faculty of Science, The University of Hong Kong, Hong Kong SAR, China 4. State Key Laboratory of Brain and Cognitive Sciences, The University of Hong Kong, Hong Kong SAR, China

INTRODUCTION

It has been observed that for complex phenotypes, associated SNPs become detectable as the sample size reaches a certain minimum, above which the number of significantly associated SNPs increases with increasing sample size in an approximately linear fashion [1]. However, a detailed modelling of the entire relationship between the number of significant SNPs and sample size has not been performed. Here, we consider a point-normal polygenic model, to predict the number of independent significant SNPs given genome-wide association study (GWAS) sample size, key parameters of genetic architecture of the phenotype and disease. Under the same assumptions, we also predicted several important outcome indices of the GWAS, including the overall phenotypic variance explained by significant SNPs, and the predictive accuracy of polygenic scores that weight SNPs by the effect size estimates from GWAS, after shrinkage by various methods. Calculations were performed from analytically derived formulae, and were validated by simulations. We compared the predictions of our method to the observed results of GWAS on height and BMI, and found that they were in agreement. Thus, our method could be a useful tool for the design of GWAS and for predicting the future behavior of GWAS as sample sizes increase further.

MODEL DESCRIPTION

1. Effect size distribution

The overall effect size across all SNPs follows a point-normal distribution.

$$\begin{cases} \beta \sim N\left(0, \frac{h^2}{m(1-\pi_0)}\right) & \text{proportion of } \pi_0 \\ \beta \sim N\left(\frac{h^2}{m(1-\pi_0)}, \frac{h^2}{m(1-\pi_0)}\right) & \text{proportion of } 1-\pi_0 \end{cases}$$

- h^2 is the SNP heritability or the heritability on the liability scale if the trait is binary.
- m is total number of nearly independent SNPs that may contribute to the phenotypic variance.
- π_0 is the proportion of SNPs that do not contributing to the variance of phenotype.

Accordingly, the distribution of effect size estimates across all SNPs is a nor

$$\begin{cases} \hat{\beta} \sim N\left(0, \frac{h^2}{m(1-\pi_0)}\right) & \text{proportion of } \pi_0 \\ \hat{\beta} \sim N\left(\frac{h^2}{m(1-\pi_0)}, \frac{h^2}{m(1-\pi_0)}\right) & \text{proportion of } 1-\pi_0 \end{cases}$$

n is sample size. When phenotype is binary, we approximate the log odds ratio is calculated by logistic regression model to effect size on the liability scale [2]. K is disease prevalence.

2. Polygenic model

To predict genetic risk, we adopt the polygenic model $y_i = \sum_{j=1}^m \beta_j x_{ij} + \varepsilon_i$. Under this model, the true and estimated polygenic scores are $G_i = \sum_{j=1}^m \beta_j x_{ij}$ and $\hat{G}_i = \sum_{j=1}^m \hat{\beta}_j x_{ij}$. $\text{corr}^2(G_i, \hat{G}_i)$ is used to measure the prediction accuracy.

METHODS

- To obtain the expectation and variance of statistical power, we first partitioned a wide range of probable effect sizes into a large number of small, equal and non-overlapping intervals.

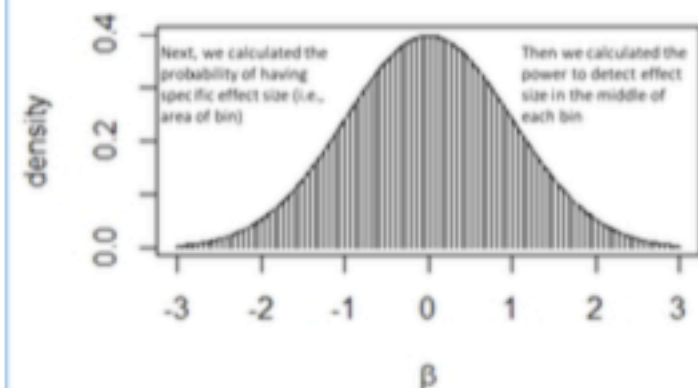


Fig. 1 Calculation of the expectation of statistical power: a demonstration

- Based on the definition of variance, we calculated the variance of power.
- The statistical powers of the intervals, weighted by their probabilities, are used to obtain a kernel density estimate of the distribution of statistical power.
- Based on the expectation and variance of statistical power, we derived the formula to calculate the expected number of significant SNPs, $E(X)$, and variance explained by the significant SNPs, $\text{Var}(X)$.
- We validated the derived formulae by simulating 100 GWASs with given parameters, counting the number of significant SNPs and calculating the variance explained by these SNPs. Comparison between the theoretical (orange and green dots) and empirical (violin plots) results is shown in Fig. 2.
- We also predict the polygenic score (PGS) prediction accuracy based on derived formulae.

APPLICATION

We compared results predicted by our model to the reported GWASs on height, BMI, and major depressive disorder and found agreement.

Phenotype	Study sample size	Number of significant SNPs		Variance explained by significant SNPs	
		Predicted	Reported	Predicted	Reported
Height	~700,000	5194	5190	14.05%	15.74(2%)
BMI	~700,000	574	716	8.5%	5.9(2%)
Major depressive disorder	~800,000	300	300	3.05%	3.5(2.2%)

Table 1. Theoretical versus reported number of independent SNPs and variance explained by these SNPs

REFERENCES

1. Wray, N. M., et al. (2017). 10 years of GWAS discovery: Biology, function, and translation. *Am J Hum Genet*, 101(1), 5-25. doi:10.1016/j.ajhg.2017.06.001. 2. Wu, T., & Sham, P. C. (2018). On the number of genetic effect loci in the light of liability scale models. *J. Stat. Theory Pract.*, 18(1), 1-15. doi:10.1080/15337817.2018.1481111. 3. Liu, Z., et al. (2018). Local Four-Category Test Height and Polygenic Score on Long COVID Score by Data Annot. *Genet. J.*, 18(1), 1-15. doi:10.1080/15337817.2018.1481111. 4. Zhang, Y. D., et al. (2018). Regression shrinkage and selection via the Lasso. *Statist. Sinica*, 18(2), 103-130. doi:10.1007/s11464-018-0671-8. 5. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8. 6. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8. 7. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8. 8. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8. 9. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8. 10. Tibshirani, R. (1996). *Regression shrinkage and selection via the Lasso*. *Statist. Sinica*, 16(1), 1-40. doi:10.1007/s11464-018-0671-8.

RESULTS

1. Distribution of statistical power to detect causal SNPs becomes a bi-model with the increase of sample size.

2. To Predict the expectation and variance of the number of independent significant SNPs

- $E(X) = m\pi_0\alpha + m(1-\pi_0)E(p)$
- $\text{Var}(X) = m\pi_0\alpha(1-\alpha) + m(1-\pi_0)[E(p) + 1 - E(p)] - \text{var}(p)$, p is statistical power.

Theoretical values or their 95% confidence intervals overlap with the empirical results, which validated our derived formulae (Fig. 2).

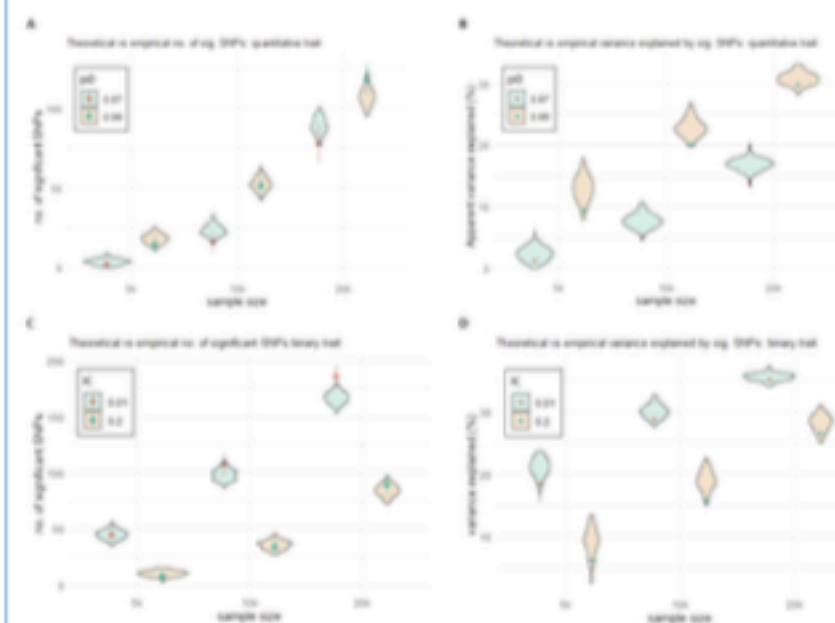


Fig. 2 Theoretical versus empirical number of independent significant SNPs and variance explained by these SNPs ($m = 50k$, $\pi_0 = 0.99$, $h^2 = 0.4$; case ratio is 0.5 for binary trait)

3. To predict the polygenic score (PGS) prediction accuracy

- Among the methods to construct PGS, the conditional expectation $E(\hat{G}_i|G_i)$ or the local TDR [3] method, is always the best in terms of the prediction accuracy.
- With moderate sample size, LASSO [4] outperforms p-value thresholding [5] and ordinary least square estimate. However, when sample size is large enough, p-value thresholding is also able to give accurate polygenic prediction.

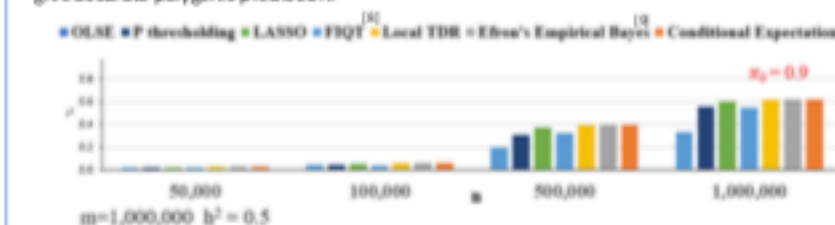


Fig. 3 Predicting prediction accuracy of PGS constructed by various of shrinkage method