

MegaPRS in Math

Tuesday, January 23, 2024 3:00 PM

[Improved genetic prediction of complex traits from individual-level data or summary statistics | Nature Communications](#)

Method

Suppose n individuals and m snps, X (size $n \times m$), X_j column contains genotypes for SNP j , Y denotes phenotype.

X_j and Y been standardized, $N(0,1)$, the error term is also $N(0,1)$

Consider the prediction of X on Y , by linear regression

$$E[Y] = X_1\beta_1 + X_2\beta_2 + \dots + X_m\beta_m = X\beta \quad (1)$$

β is the effect size of SNP j , and since X_j and Y are standardized, so we have $h_j^2 = \beta_j^2$

Heritability

$$h^2 = \frac{Var(X\beta)}{Var(Y)}$$

The heritability model takes form:

$$E[h_j^2] = a_{j1}\tau_1 + a_{j2}\tau_2 + \dots + a_{jK}\tau_K \quad (2)$$

where the a_{jk} are pre-specified SNP annotations, while the parameters τ_k are estimated from the data.

So if we take $h^2 = \tau_1 I_j (f_j(1-f_j))^{0.75}$, then the variance of β_j is $\tau_1 I_j (f_j(1-f_j))^{0.75} \times h^2$, by SumHer

$$\beta_j \sim N(0, \tau_1 I_j (f_j(1-f_j))^{0.75} \times h^2)$$

$$\hat{Y} = \sqrt{h^2} X\beta + \sqrt{1-h^2} \epsilon$$

Calculation of heritability

Suppose I have X and Y as individual data, we can calculate h^2 by the following formulas:

1. $r = \frac{X^T Y}{n}$, this is the correlation between X and Y
2. $S = n \frac{r^2}{1-r^2}$, this is the chi-square test statistics for SNP j .
3. For SumHer, it calculates τ by the following formula:

$$E[S_j] \approx 1 + n \sum_l c_{jl}^2 (a_{l1}\tau_1 + a_{l2}\tau_2 + \dots + a_{lK}\tau_K) \quad (3)$$

- a. S was calculated by step 2. C_{jl}^2 is taken from the correlation matrix of SNPs j and l .
 - b. We can take a linear regression of S and Correlation, and the coefficient is τ
4. $h^2 = \tau I_j (f(1-f))^{0.75}$, Where I is the LD matrix, and f is the MAF.

Suppose we are using Summary Statistics instead of individual data, the first step will be updated:
We have genotype data and Z-score.

$$S = Z^2$$

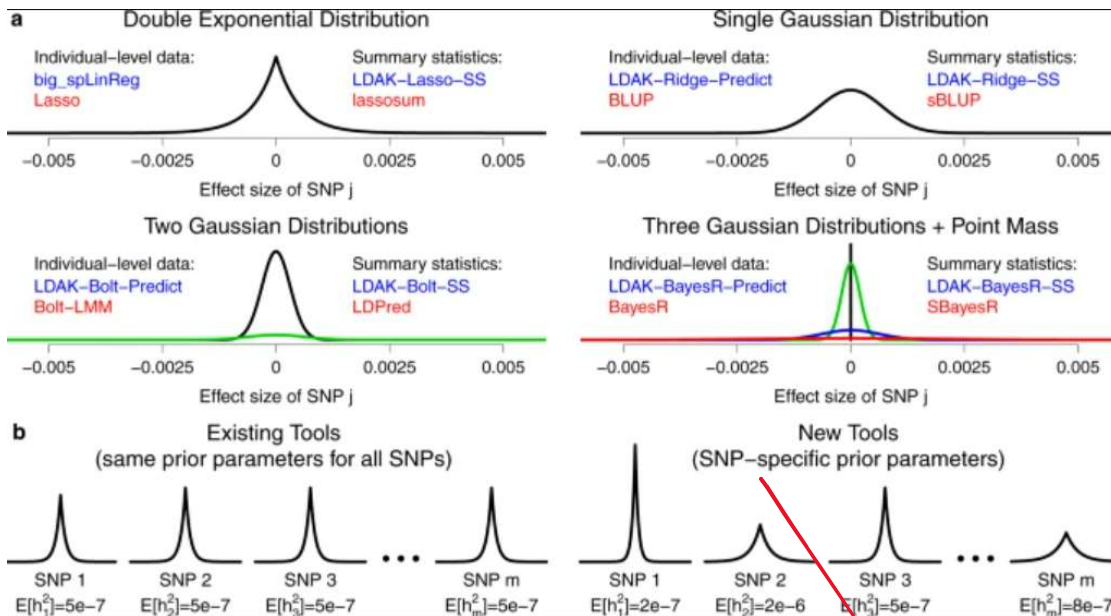
we can skip calculating the correlation between X and Y and directly get $S = Z^2$.

Prediction Model for β

Y in fact is based on $X\beta$ and ϵ , and ϵ is normally distributed.

$$Y \sim N(X\beta, \sigma_\epsilon^2)$$

We have h^2 already, so we need to estimate β . The newest methods estimate β by using Bayes Inference. Different methods would have different prior distributions.



1. Double Exponential Distribution

$$\beta_j \sim \text{DE}(\lambda / E[h_j^2]^{0.5})$$

2. Single Gaussian Distribution (Ridge Regression)

$$\beta_j \sim N(0, E[h_j^2]) \text{ and } \beta_j \sim N(0, \nu E[h_j^2])$$

Respectively for using Individual data or Summary Statistics.

3. Two Gaussian Distribution (Bolt-lmm)

Predict and LDAK-Bolt-SS use a mixture of two Gaussian distributions, $\beta_j \sim p N(0, (1-f_2)/p E[h_j^2]) + (1-p) N(0, f_2/(1-p) E[h_j^2])$.

4. Three Gaussian Distributions + Point Mass

LDAK-BayesR-Predict and LDAK-BayesR-SS use a mixture of a point mass at zero and three Gaussian distributions, $\beta_j \sim \pi_1 \delta_0 + \pi_2 N(0, sE[h_j^2]/100) + \pi_3 N(0, sE[h_j^2]/10) + \pi_4 N(0, sE[h_j^2])$, where $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$ and $s = (\pi_2/100 + \pi_3/10 + \pi_4)^{-1}$.

Calculation for β

Assuming we have $\hat{\beta}$ from Summary Statistics ($\beta = nZ$) as $P(\hat{\beta})$, and we get the prior distribution of β as $P(\beta)$

For Bayesian Inference:

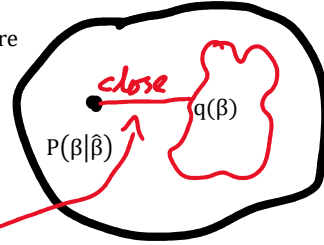
$$P(\beta|\hat{\beta}) = \frac{P(\hat{\beta}|\beta)P(\beta)}{P(\hat{\beta})}$$

$$\text{Posterior} = \frac{\text{MLE} \times \text{Prior}}{\text{Observe}}$$

So we need to calculate the MLE of β and $\hat{\beta}$, for example by Gradient Descent MLE, while these methods are slow in general. And posterior has no closed form, so it's not feasible to compute $P(\beta|\hat{\beta})$ directly.

Posterior — MLE Prior
Observe

So we need to calculate the MLE of β and $\hat{\beta}$, for example by Gradient Descent MLE, while these methods are slow in general. And posterior has no closed form, so it's not feasible to compute $P(\beta|\hat{\beta})$ directly.



Using Variational Bayes

Our goal becomes to find a distribution $q(\beta)$ that approximates the posterior distribution $P(\beta|\hat{\beta})$

■ Variational Bayes (VB): L is KL divergence

$$L(q(z), p(z|x)) = \text{KL}(q(z)||p(z|x))$$

$$\text{KL}(q(\beta)||P(\beta|\hat{\beta})) = \int q(\beta) \log \frac{q(\beta)}{P(\beta|\hat{\beta})} d\beta$$

$$= \int q(\beta) \log \frac{q(\beta)P(\hat{\beta})}{P(\beta, \hat{\beta})} d\beta$$

$$= - \int q(\beta) \log \frac{P(\beta, \hat{\beta})}{q(\beta)} d\beta + \log P(\hat{\beta})$$

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

The first part is called ELBO (evidence lower bound), and since we need to minimize KL divergence, we can drop the second term $\log P(\hat{\beta})$.

$$q(\beta) = \text{argmax}(\text{ELBO}(q(\beta)))$$

And $q(\beta)$ is the distribution of β that we need.

Prediction

We have β and h^2 , and we have the genotype data, we can then calculate PRS score:

$$\hat{Y} = \sqrt{h^2}X\beta + \sqrt{1-h^2}\epsilon$$

Discussion

About β

If we standardize X and Y into $N(0,1)$, then β will be $N(0,1)$

$$Z = \beta / SE(\beta) = \sqrt{n}\beta \text{ and } SE(\beta) = 1/\sqrt{n}$$

Because $SE = SD/\sqrt{n}$

Steps of MegaPRS

Correlation Matrix and High LD matrix

Date. 2024/1/28

MegaPRL each step, with input and output.

- calculation of high-lol and correlation matrix.

1. ① Correlation Matrix.

Input, randomly pick ~~two~~ ^{individuals} with bin/bed/fam.

n	ind 1	0	1	2	0	1	2
	ind 2	0	1	2	1	0	2

	ind two	0	0	1	2	1	0

② Calculation: Pearson's correlation's coefficient.

• SNP_i and SNP_j .

$$r_{ij} = \frac{\sum_{k=1}^n (X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j)}{\sqrt{\sum_{k=1}^n (X_{ki} - \bar{X}_i)^2 \sum_{k=1}^n (X_{kj} - \bar{X}_j)^2}}$$

③ Output: $m \times m$ matrix.

II. ① High LD matrix.

Input: reference panel of N inds and M snps.

High LD region list.

Index	From	End	chr.
region 1	47888888	52111111	1
...
region 24	-	-	-

- We'll compute high LD of SNPs in these regions.

② Calculation.

• P_{AB} , P_{AB} , P_{AB} , P_{AB} .

- D'_{ij} represents LD of SNP_i and SNP_j .

$$\checkmark D'_{ij} = \frac{P_{AB} P_{ab} - P_A P_a P_B P_b}{\min[P_A P_a, P_B P_b]} \quad \text{--- ChatGPT}$$

$[-1, 1]$

$$\checkmark D^+_{ij} = P_{AB} - P_A P_B \quad \text{--- PopGen Note}$$

$$\checkmark r^2 = \frac{D^2}{P_A(1-P_A)P_B(1-P_B)} \quad \text{--- PopGen Note}$$

$[0, 1]$

③ Output: $M \times M$ matrix.

二. Prediction Model

① Input: summary statistics

cor matrix

LD matrix

→ model bayesT.

② Predicting by BayesT, h^2 .

$$h^2 = \sum_i I_{ij} (f_{ij}(1-f_{ij}))^{0.75}.$$

$$I_{ij} = E I_{ij} / N, \quad f_{ij} = E f_{ij} / N.$$

and τ_i is set to be 0.1, 0.3, 0.5.

$$\Rightarrow h_1^2, h_2^2, h_3^2.$$

③ Predicting by BayesT, β .

$$\text{Prior: } \beta_j \sim \tau_1 \delta_0 + \tau_2 N(0, s[h_j^2]/w) + \tau_3 N(0, s[h_j^2]/10) + \tau_4 N(0, s[h_j^2]).$$

$$\tau_1 + \tau_2 + \tau_3 + \tau_4 = 1, \quad s = (\tau_2/w + \tau_3/10 + \tau_4)^{-1}.$$

and τ_2, τ_3, τ_4 take from (0, 0.01, 0.05, 0.1, 0.2)

can be replicated. \Rightarrow 35 models.

$$\Rightarrow 35 \times 3 = 105 \text{ models}.$$

③ Posterior $P(\beta | \hat{\beta})$.

Assuming we have $\hat{\beta}$ from summary statistics (standardised),
 $P(\hat{\beta})$ represents it ($\delta=1$, $\hat{\beta} = \frac{1}{\sqrt{n}} z$)
 and $P(\beta)$ as prior distribution from ②.

$$P(\beta | \hat{\beta}) = \frac{P(\hat{\beta} | \beta) P(\beta)}{P(\hat{\beta})}$$

⇒ Using VB to approximate $P(\beta | \hat{\beta})$.

Assuming $q(\beta)$ approximates $P(\beta | \hat{\beta})$.

$$q(\beta) = \operatorname{argmax} \left(+ \int q(\beta) \log \frac{P(\beta, \hat{\beta})}{q(\beta)} d\beta \right)$$

三. ④ Calculating PRL.

Input, Model and genotype data.

$$\hat{Y} = \sqrt{h^2} X \beta + \sqrt{1-h^2} \varepsilon \quad \varepsilon \sim N(0, 1)$$

四. Comparison: correlation (Y, \hat{Y}) .