



JOHNS HOPKINS
UNIVERSITY

An ensemble penalized regression method for multi-ancestry polygenic risk prediction

Jingning Zhang*, Jianan Zhan, Jin Jin, Cheng Ma, Ruzhang Zhao, Jared O'Connell,

Yunxuan Jiang, 23andMe Research Team, Bertram L. Koelsch, Haoyu Zhang, **Nilanjan Chatterjee***

*Correspondence to: **Jingning Zhang (jzhan218@jhu.edu)** and **Nilanjan Chatterjee (nilanjan@jhu.edu)**

Penalized regression for PRS construction

- ▶ We propose to solve the PRS coefficients in a penalized regression
- ▶ If we have individual-level data, estimate of β can be obtained by minimizing

$$\frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \text{penalty}(\beta)$$

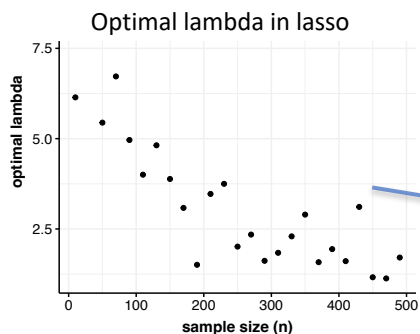
- ▶ However, we usually don't have individual-level data
 - ▶ By approximating $\mathbf{X}^T \mathbf{X}$ by LD, \mathbf{D} , in a separate reference data, and $\mathbf{X}^T \mathbf{y}$ by summary statistics, $\hat{\beta}$, it is equivalent to minimizing

$$\beta^T \mathbf{D} \beta - 2\beta^T \hat{\beta} + \text{penalty}(\beta)$$

Polygenic Risk scOres based on enSemble of PEnalized Regression models (PROSPER)

- Solution of PRS coefficient for K populations, $\beta_i, i = 1, \dots, K$, can be obtained using **coordinate decent algorithm** by minimizing

$$\sum_{i=1}^K \beta_i^T \underset{\text{LD matrix}}{\mathbf{D}_i} \beta_i - 2 \beta_i^T \underset{\text{Summary statistics}}{\hat{\beta}_i} + \sum_{i=1}^K \underset{\substack{\text{Tuning parameter} \\ \text{for sparsity}}}{\lambda_i} |\beta_i| + \sum_{1 \leq i_1 < i_2 \leq K} \underset{\substack{\text{Tuning parameter for} \\ \text{genetic similarity}}}{c_{i_1 i_2}} (\beta_{i_1} - \beta_{i_2})^2$$



- Tuning parameter

- $\lambda_i, i = 1, \dots, K$

- $c_{i_1 i_2}, i_1 < i_2 = 1, \dots, K$

- **How to reduce the number of grid search when K is large?**

- Let $\lambda_i = \lambda \lambda_{i_0}$
where λ_{i_0} is the value used in lasso (single-ancestry)

- Let $c_{i_1 i_2} = c$

- **Super learning**

- Weighted combine PRS generated from all tuning parameter settings for all ancestries

- # PRS combined

$(\# \text{ candidate values})^{(\# \text{ tuning parameters})} \times (\# \text{ ancestries})$

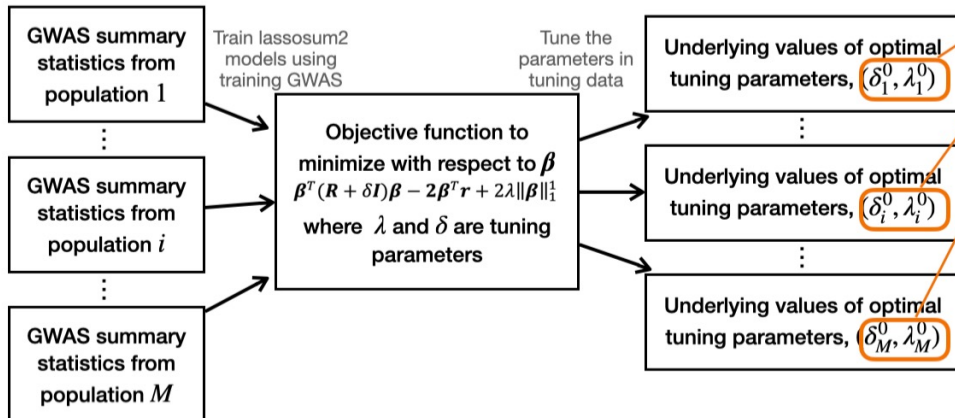
E.g. 5-ancestry analysis

- 2 tuning parameters: λ and c ; 5 candidate values each

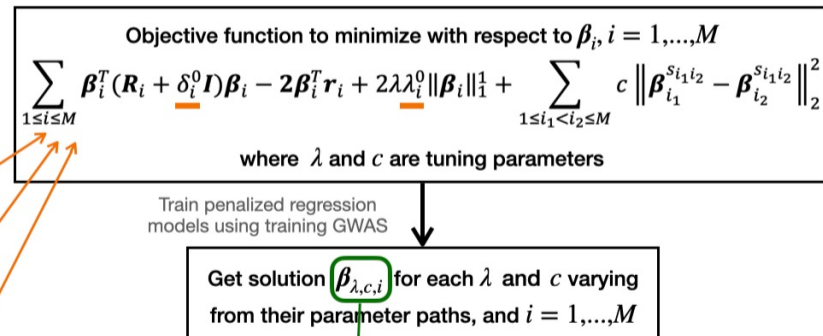
- # PRS combined by super learning: $5^2 \times 5 = 125$

Flowchart

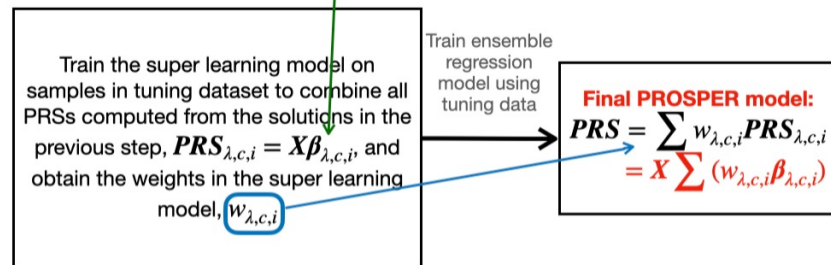
Step 1: Separate single-ancestry analysis for all populations



Step 2: Joint analysis across populations using penalized regression

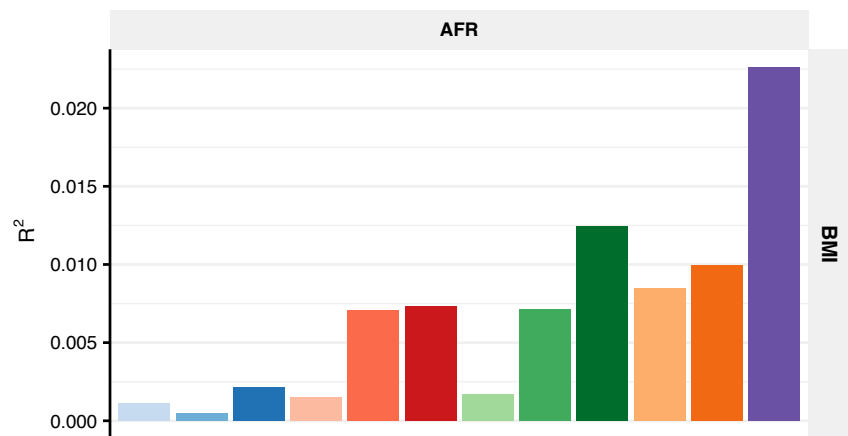


Step 3: Ensemble regression

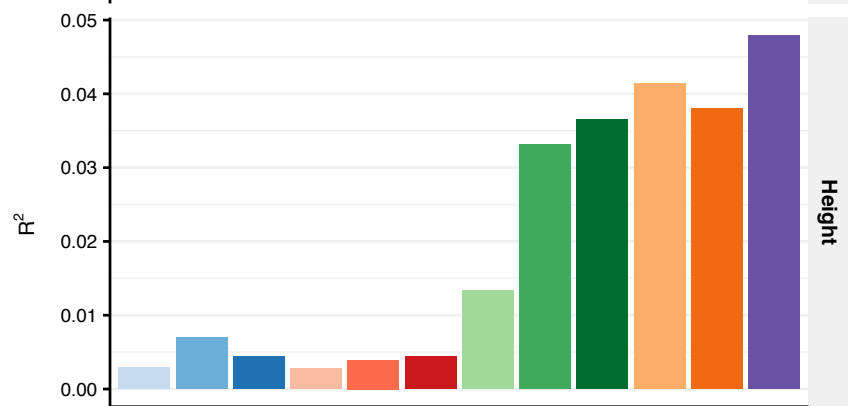


Results on data from All of US (AoU)

a



b



Single ethnic method

- CT
- LDpred2
- lassosum2

EUR PRS based method

- EUR CT
- EUR LDpred2
- EUR lassosum2

Multi ethnic method (weighted PRS)

- weighted CT
- weighted LDpred2
- weighted lassosum2

Multi ethnic method (existing methods)

- PRS-CSx
- CT-SLEB

PROSPER

- PROSPER

- Averaged sample size
 - EUR (N = 48K)
 - AFR (N = 22K)

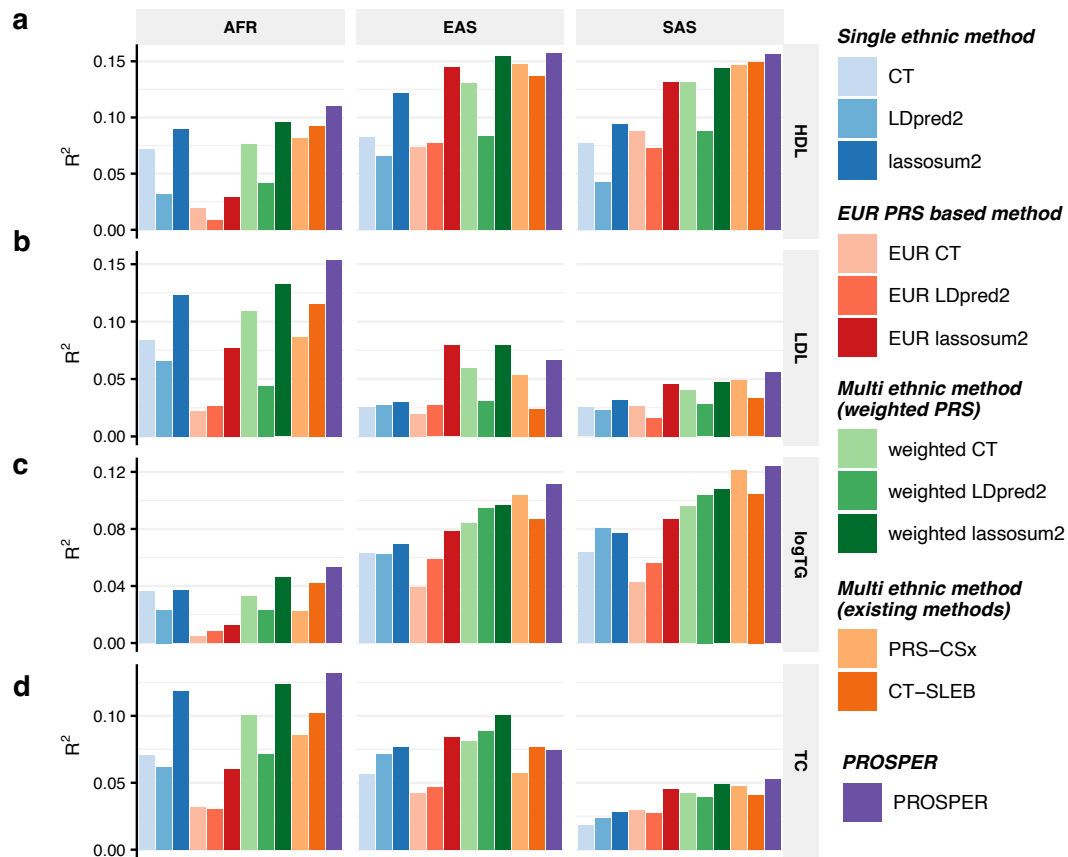
- Improvement PROSPER over alternative methods (average across traits and ancestries) in R^2
 - 91.3% over PRS-CSx
 - 76.5% over CT-SLEB
 - 56.8% over weighted lassosum2

SNP set:

HM3 for PRS-CSx

HM3+MEGA for all other methods

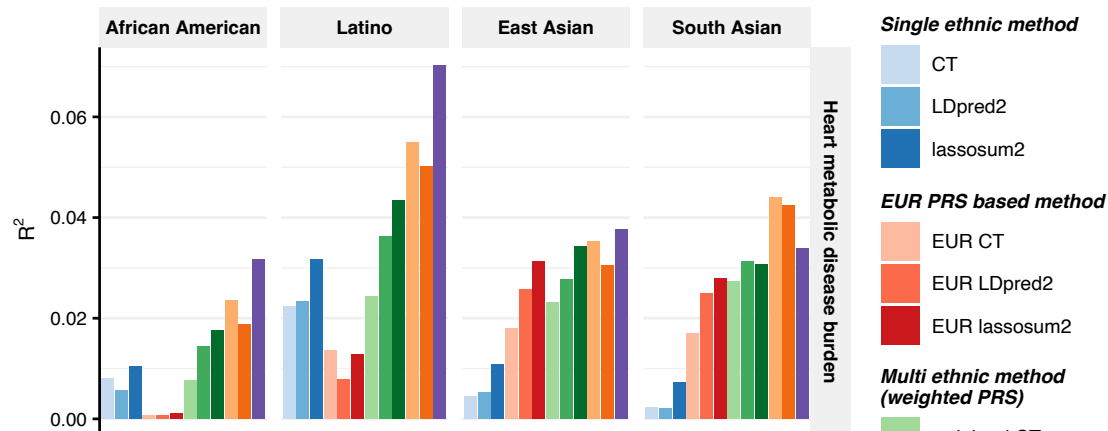
Results on data from Global Lipids Genetics Consortium (GLGC)



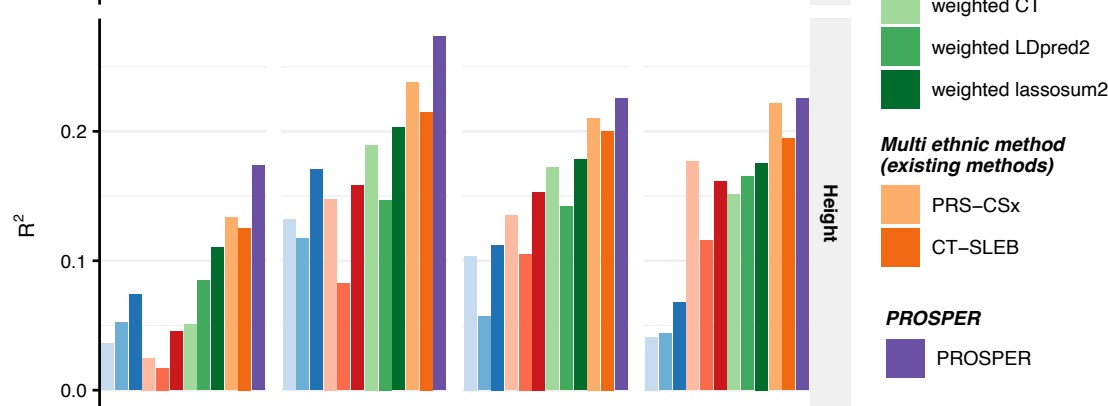
- Averaged sample size
 - EUR (N = 931K)
 - AFR (N = 93K)
 - EAS (N = 146K)
 - SAS (N = 34K)
- Improvement of PROSPER over alternative methods (average across traits and ancestries) in R^2
 - 34.2% over PRS-CSx
 - 37.7% over CT-SLEB
 - 13.1% over weighted lassosum2 in AFR
 - 12.3% over weighted lassosum2 in SAS

Results on data from 23andMe Inc. (23andMe) continuous traits

a



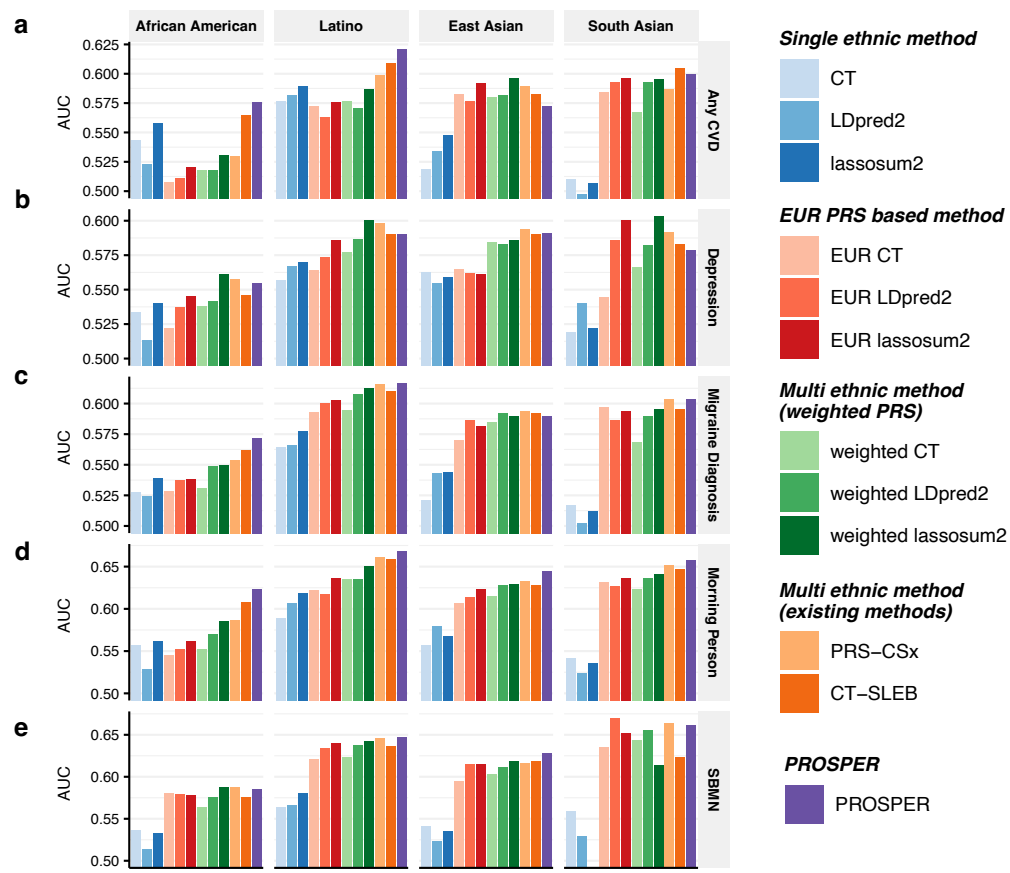
b



- ▶ Averaged sample size
 - ▶ EUR (N = 2700K)
 - ▶ African American (N = 136K)
 - ▶ Latino (N = 442K)
 - ▶ East Asian (N = 116K)
 - ▶ South Asian (N = 31K)

- ▶ Improvement of PROSPER over alternative methods (average across traits) in R^2
 - ▶ 32.4% over PRS-CSx in African American
 - ▶ 21.5% over PRS-CSx in Latino

Results on data from 23andMe Inc. (23andMe) binary traits



- ▶ Averaged sample size
 - ▶ EUR (N = 2370K)
 - ▶ African American (N = 109K)
 - ▶ Latino (N = 401K)
 - ▶ East Asian (N = 86K)
 - ▶ South Asian (N = 24K)
- ▶ Improvement of PROSPER over alternative methods (average across traits and ancestries) in logit-scale variance
 - ▶ 12.3% over PRS-CSx
 - ▶ 7.8% over CT-SLEB

Computational and memory usage

- ▶ Chromosome 22
- ▶ two ancestries: AFR, and EUR
five ancestries: AFR, AMR, EAS, EUR, and SAS
- ▶ Computational time and memory usage

Method	Computational time (minutes)	Memory (Gb)
PROSPER (two ancestries)	3.0	2.24
PROSPER (five ancestries)	6.8	2.35
PRS-CSx (two ancestries)	111.1	0.78
PRS-CSx (five ancestries)	595.8	0.84

Conclusion

- ▶ PROSPER has substantial improvement over alternative methods and can deal with traits with a variety of genetic architectures
- ▶ PROSPER is not only more powerful for complex traits **with high polygenicity**, but also **robust to biomarker traits with large-effect loci**. In addition, PROSPER developed based on penalized regression is **an order of magnitude faster** compared to alternative Bayesian methods.