

A Fast and Scalable Framework for Large-scale and Ultrahigh-dimensional Sparse Regression with Application to the UK Biobank

<https://www.biorxiv.org/content/10.1101/630079v2.full.pdf>

INTRODUCTION

PREMISE FOR PAPER

- ▶ Biobanks = huge data sources. Massive potential.
BUT... **ultra-high dimensional**
- ▶ Shape of data drives modelling decision...
- ▶ “BIG”

Var_1	Var_2	Var_3	Var_4	Var_5	Var_6	...	Var_50000
-------	-------	-------	-------	-------	-------	-----	-----------

Feature selection
e.g. Lasso

- ▶ “BIG”

Machine Learning

Samp_1

Samp_2

Samp_3

Samp_4

Samp_5

Samp_6

...

Samp_5000000

INTRODUCTION

PREMISE FOR PAPER

- ▶ Biobanks = huge data sources. Massive potential.
BUT... **ultra-high dimensional**
- ▶ Shape of data drives modelling decision...

▶ “BIG”

▶ “BIG”

Var_1 Var_2 Var_3 Var_4 Var_5 Var_6 ... Var_50000

Feature selection
e.g. Lasso

Machine Learning

Samp_1

Samp_2

Samp_3

Samp_4

Samp_5

Samp_6

...

Samp_5000000

INTRODUCTION

REGULARIZATION

- Constraints applied to shrink regression coefficients

↑ Model stability ↓ Overfitting

- **LASSO**

- L1 penalty
- Shrinks **to** 0
- Variable selection
- Sum of coefficients = penalty

Tuning parameter

$$\lambda \sum_{j=1}^m |\hat{\beta}_j|$$

- **RIDGE**

- L2 penalty
- Shrinks **towards** 0
- Multicollinearity → grouping
- SS coef = forms penalty

$$\lambda \sum_{j=1}^m \hat{\beta}_j^2$$

- **ELASTIC NET**

- Alpha parameter
- $\alpha=0$: Ridge
- $\alpha=1$: Lasso
- $0 < \alpha < 1$: Combination

Tuning Ridge Lasso

$$\lambda \left(\frac{1 - \alpha}{1} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)$$

INTRODUCTION

REGULARIZATION

- ▶ Constraints applied to shrink regression coefficients
 - ↑ Model stability ↓ Overfitting

▶ LASSO

- L1 penalty
- Shrinks **to** 0
- Variable selection
- Sum of coefficients = penalty

$$\lambda \sum_{j=1}^m |\hat{\beta}_j|$$

Tuning
parameter

▶ RIDGE

- L2 penalty
- Shrinks **towards** 0
- Multicollinearity → grouping
- SS coef = forms penalty

$$\lambda \sum_{j=1}^m \hat{\beta}_j^2$$

▶ ELASTIC NET

- Alpha parameter
- $\alpha=0$: Ridge
- $\alpha=1$: Lasso
- $0 < \alpha < 1$: Combination

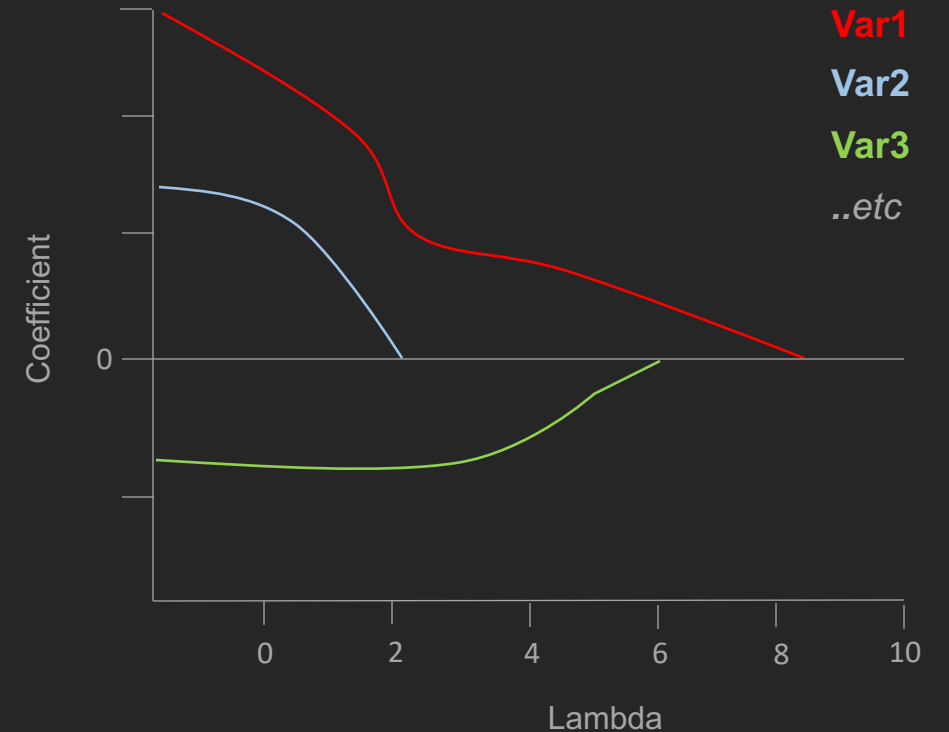
$$\lambda \left(\frac{1 - \alpha}{\sum_{j=1}^m \hat{\beta}_j^2} + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)$$

Tuning Ridge Lasso

INTRODUCTION

GLMNET

- ▶ Existing package in R, developed by Hastie and co.
- ▶ Computes Lasso coefficient paths.
(optimal lambda unknown → compute over a grid of lambda values)
- ▶ Cross validation to select optimal lambda.
- ▶ So this is not a new idea...
BUT... scaling remains an issue



INTRODUCTION

RELEVANCE TO GENOMIC PREDICTION



- ▶ Genetic variance currently explained by GWAS

RELEVANCE TO GENOMIC PREDICTION

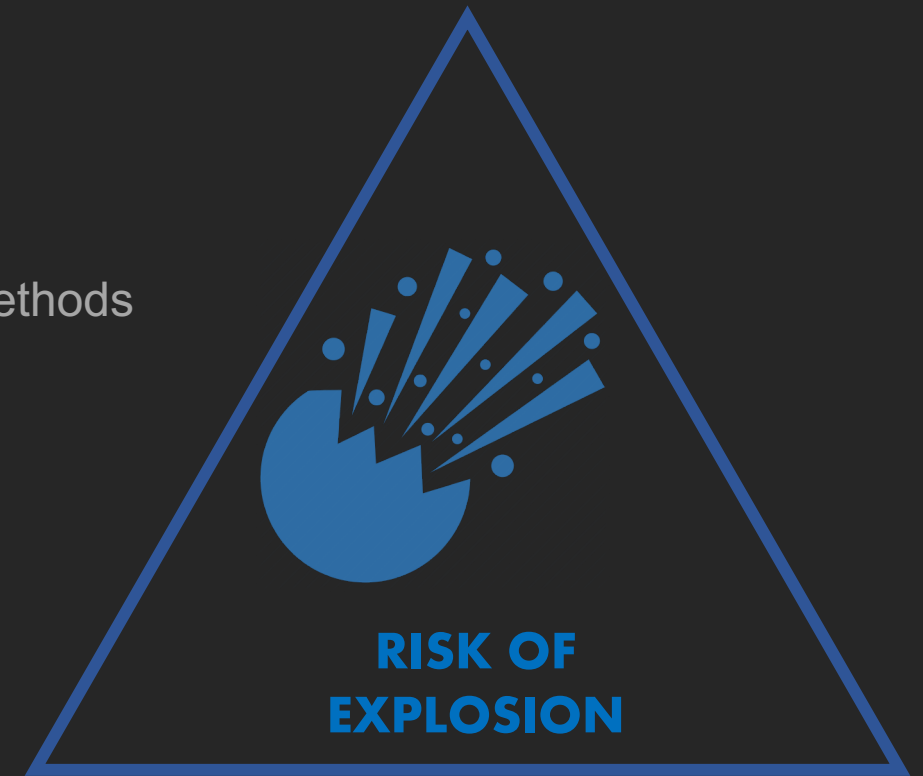


- ▶ 500,000 samples → more predictive power
- ▶ Limitations of univariate approaches

RELEVANCE TO GENOMIC PREDICTION



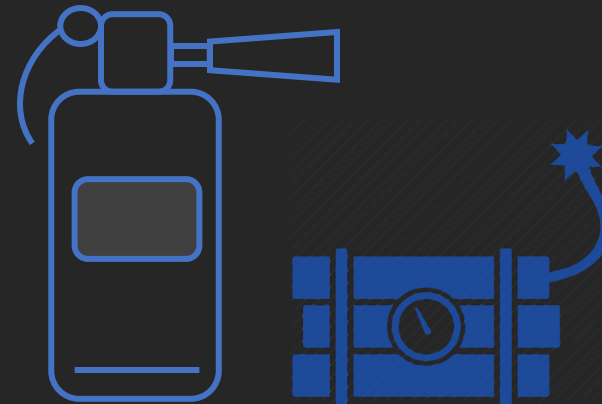
- ▶ Full-scale multivariable methods



INTRODUCTION

SNPNET

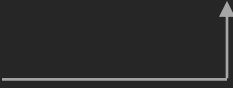
- ▶ Extension of GLMNET
- ▶ **BA**tch **S**creening **I**terative **L**asso → BASIL
- ▶ R package → easy acquisition
- ▶ Full-scale multivariable methods



METHODS

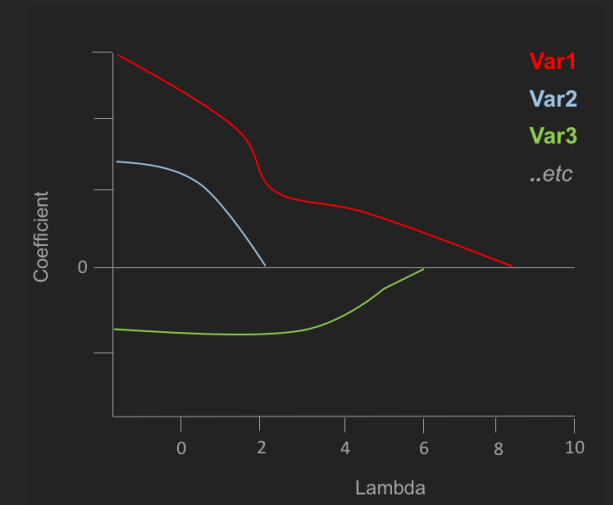
BASIL

Screen \rightarrow solve \rightarrow check

- ▶ Screen: Strong set
Inner product at $\lambda l > \lambda l_{+1} - (\lambda l - \lambda l_{+1})$
- ▶ Solve: Computed strong set
Compute lambda solution only on these subset
- ▶ Check: Entire variable matrix
Evaluate against KTT condition 🙌
- ▶ **BATCH** screening of 
multiple lambdas \rightarrow single whole data check

Not new – all done in GLMNET

Specific to SNPNET



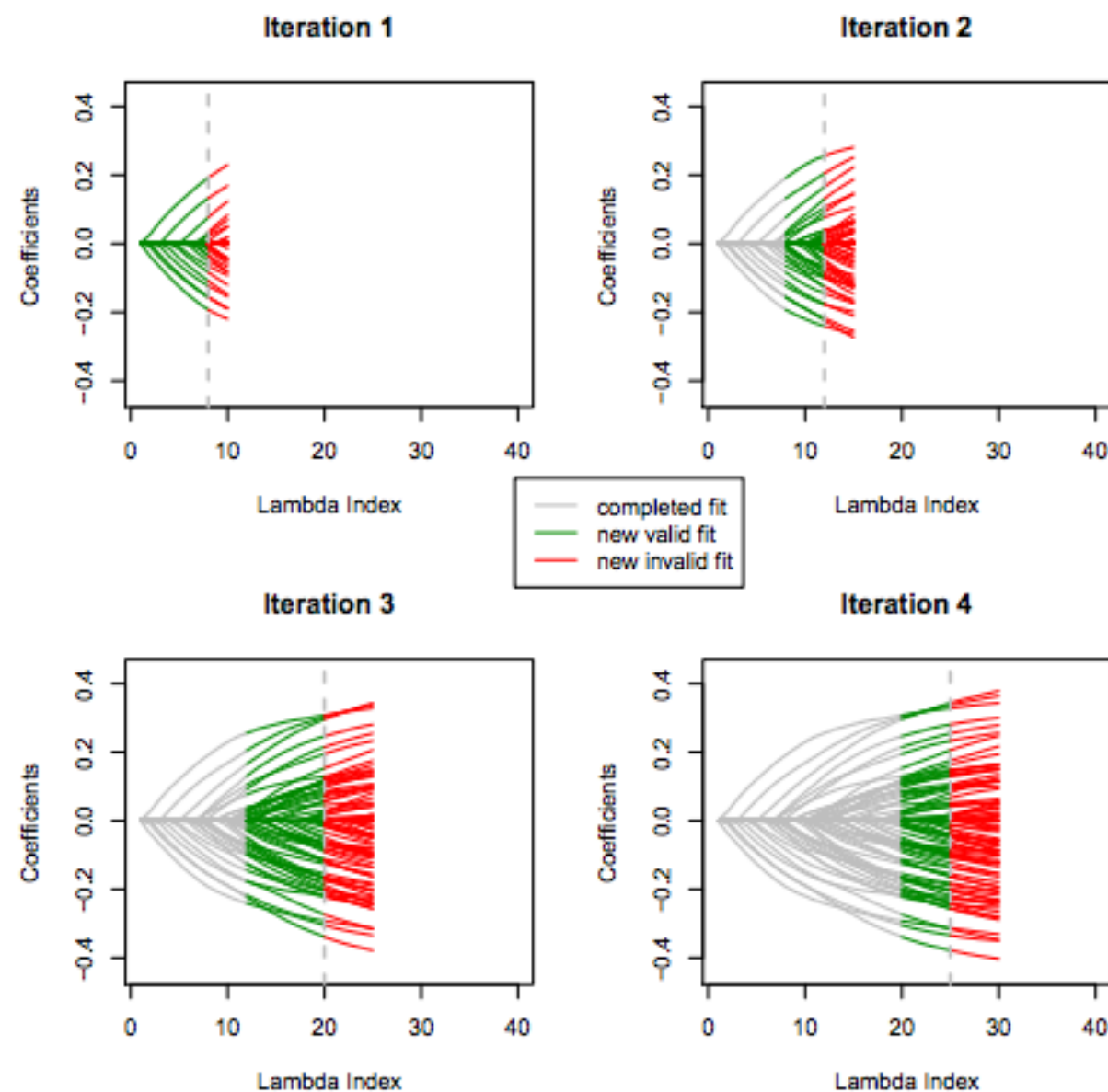


Figure 1: The lasso coefficient profile that shows the progression of the BASIL algorithm. The previously finished part of the path is colored grey, the newly completed and verified is in green, and the part that is newly computed but failed the verification is colored red.

METHODS

BASIL

- ▶ Linear
- ▶ Logistic
- ▶ Survival
- ▶ Elastic net

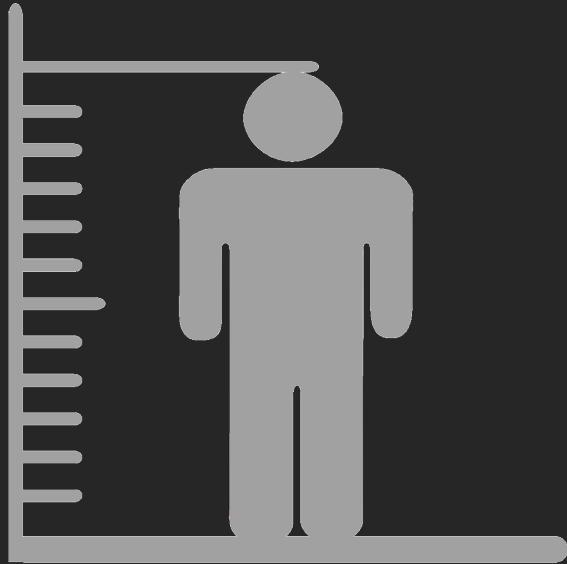
TESTING ON UKBB DATA

- ▶ 337,199 unrelated individuals (of the 500,000 in UKBB)
- ▶ Training | Validation | Test
- ▶ 805,426 measured variants
- ▶ Covariates = age, sex, principal components (40? 10?)

TESTING ON UKBB DATA

- ▶ Continuous phenotypes → R^2

1.



2.



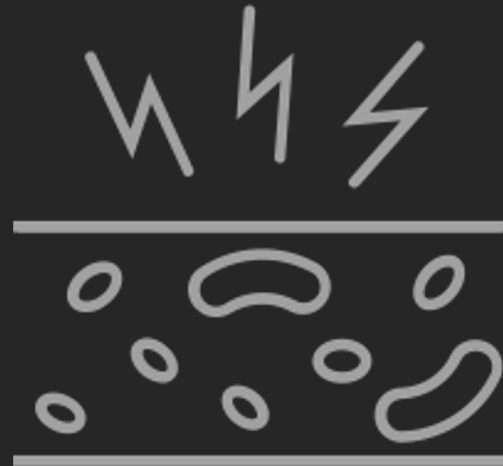
TESTING ON UKBB DATA

- ▶ Binary phenotypes → AUC


1.




2.



COMPARATIVE MODELS

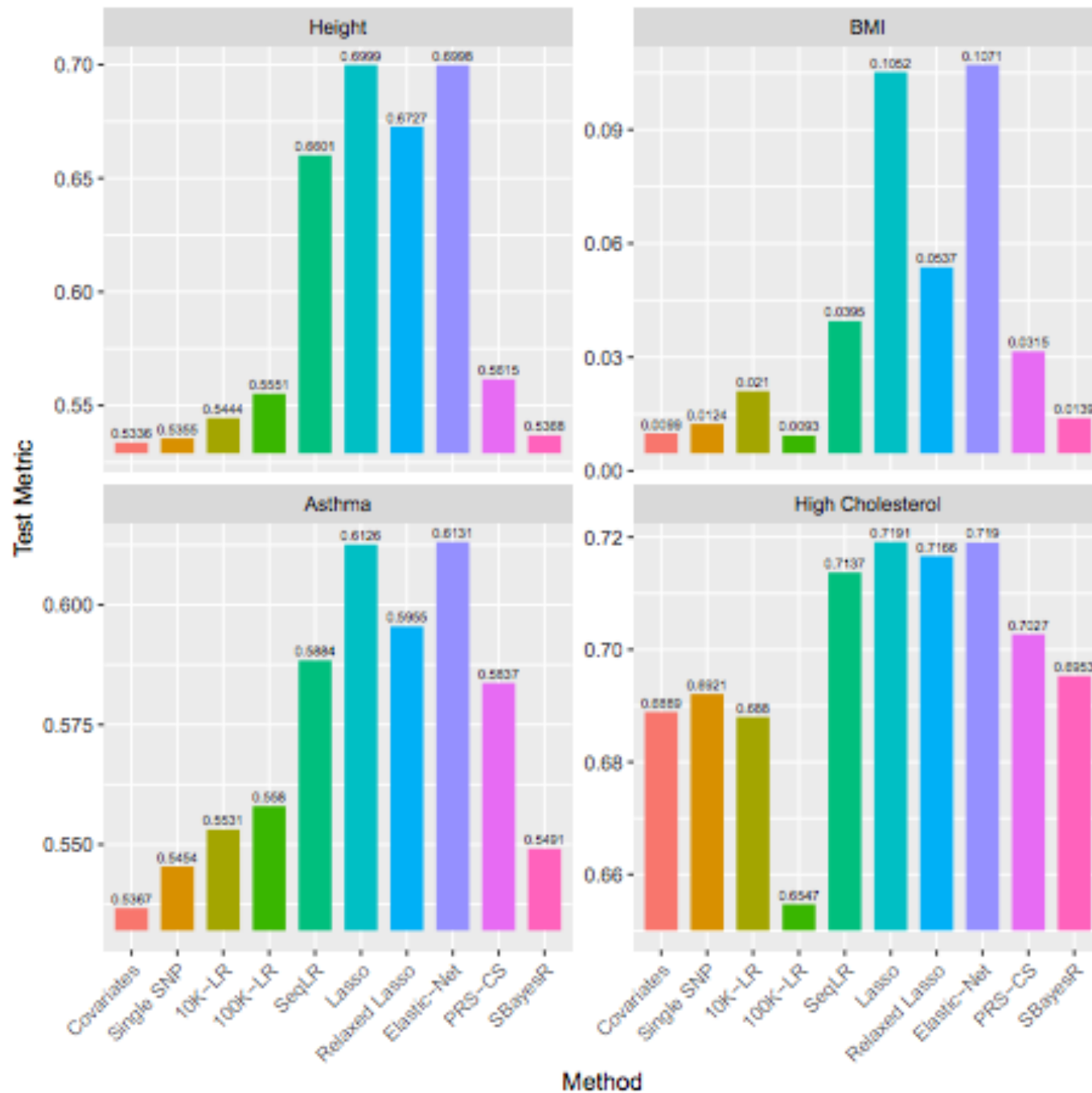
Model 1	Model 2	Model 3	Model 4	Model 5	Model 6	Model 7	Model 8	Model 9	Model 10
Baseline	Covariate only	Single SNP	Relaxed threshold 10k	Relaxed threshold 100k	Sequential	PRS-CS	SBayesR	Elastic net	Relaxed LASSO
Sex Age	Sex Age 10 PCs	Strongest univariate +model 2	Single var from linear combinations using univariate coefficients +model 2		Multiple regression - sequential increase SNP N	Default settings 		$\alpha = 0.1$ $\alpha = 0.5$ $\alpha = 0.9$	OLS variables selected for the LASSO

Most similar to PRSice 

RESULTS

EXTRACT FROM PAPER

Results on test-sample set



R2

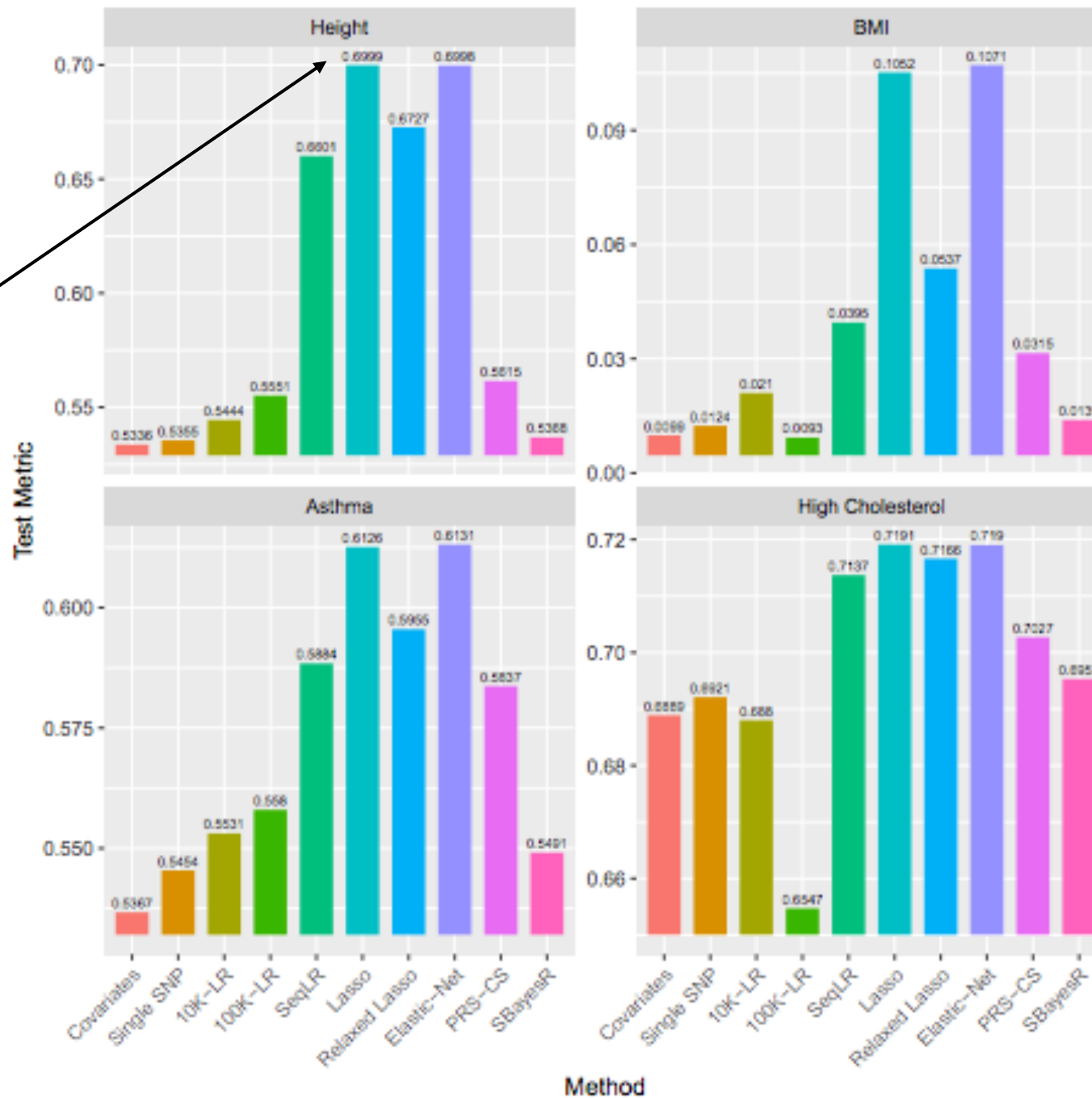
AUC

RESULTS

EXTRACT FROM PAPER

Results on test-sample set

$R^2=0.6999...$



R^2

AUC

DISCUSSION

STRENGTHS

- ▶ Thorough methodological overview
- ▶ Builds on existing method → package readily available
- ▶ Examples using real, accessible, relevant data
- ▶ Comparisons with potential alternatives
- ▶ Continuous and binary outcomes demonstrated

POTENTIAL LIMITATIONS

- ▶ Not fully clear whether individual level data required
- ▶ Examples based on split datasets...
Cross validation?
- ▶ Model fit unusually large? Covariate effect?
- ▶ Linkage disequilibrium?
- ▶ Cherry picking traits?

DISCUSSION

YOUR THOUGHTS?