



Bayesian Methods & Advanced Models for Polygenic Scores (PGS)

Jian Zeng

j.zeng@uq.edu.au



- Understanding to polygenic scores (PGS)
- Basic method to construct PGS
- Bayesian methods [15min Lecture | 15min Practical | 5min break]
- Using GWAS summary statistics [15min Lecture | 15min Practical | 5min break]
- Incorporating functional genomic annotations [15min Lecture | 15min Practical]
- Wrap-up/discussion/questions

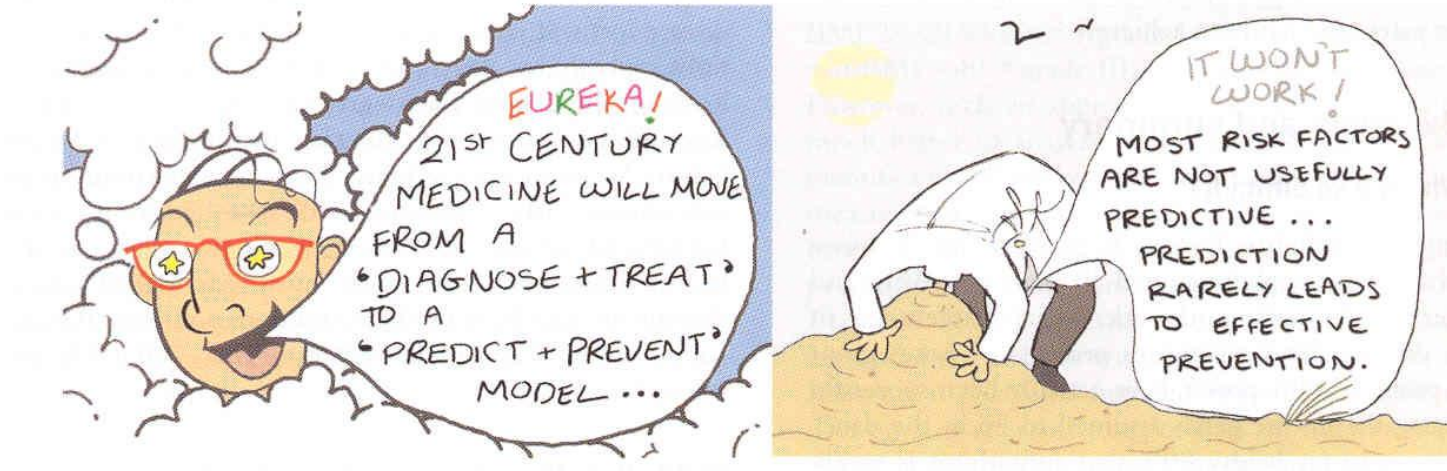
Practical exercises will use a toy example and R scripts

All materials are available at:

<https://gctbhub.cloud.edu.au/data/teaching/>

Polygenic scores (PGS) predict individual genetic values of complex traits using genome variations.

Polygenic risk scores (PRS) are predictors of the genetic susceptibilities of individuals to diseases.



Head in the clouds

Head in the sand

Source: Strachan & Read Human Molecular Genetics 3.

Theory and methodology
of polygenic scores (PGS)
are built on
our understanding of
“polygenicity”
in complex traits.



Height

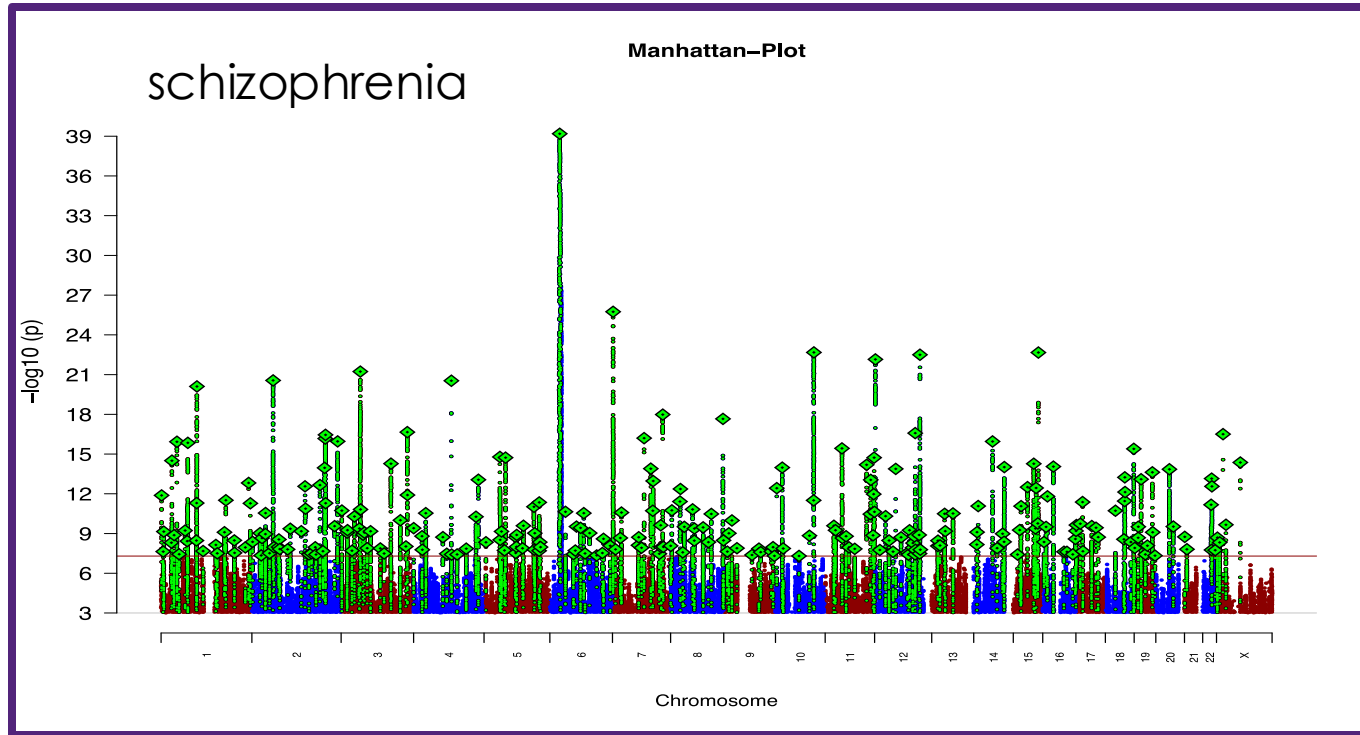


Schizophrenia



Obesity

Common diseases are polygenic



248 risk loci identified at genome-wide significance level.

We predict thousands are associated with schizophrenia.

nature

[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [articles](#) > [article](#)

Article | [Published: 08 April 2022](#)

Mapping genomic loci implicates genes and synaptic biology in schizophrenia

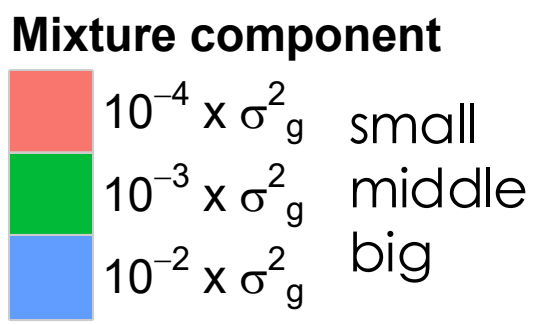
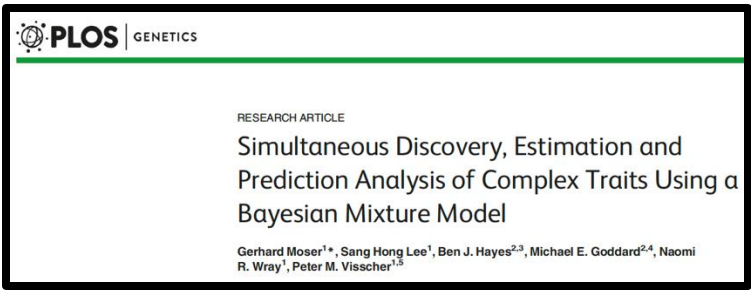
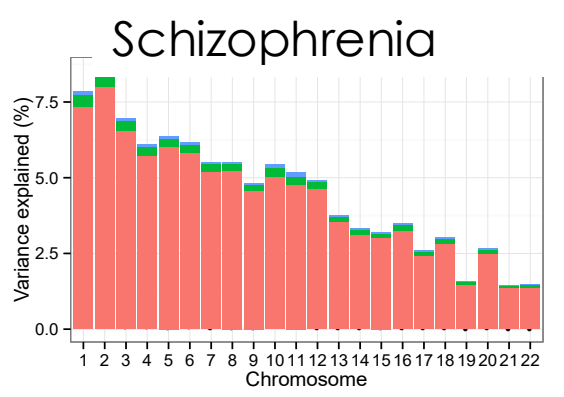
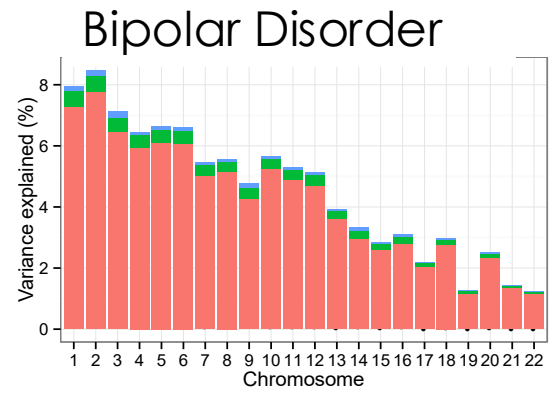
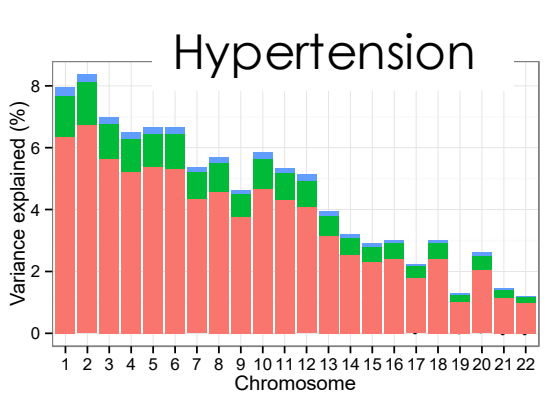
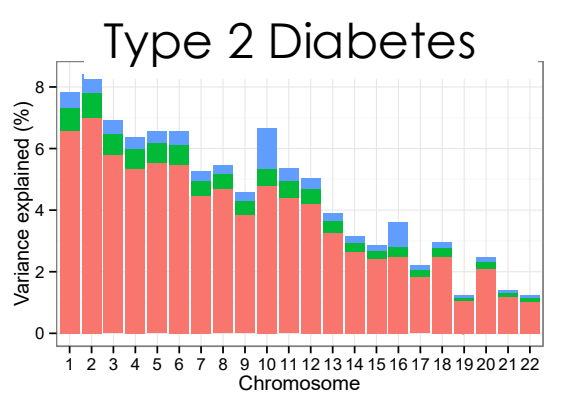
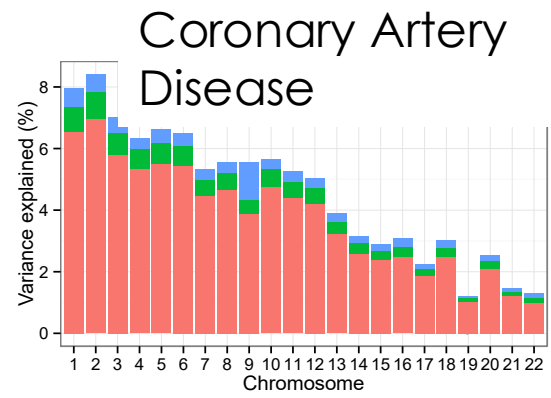
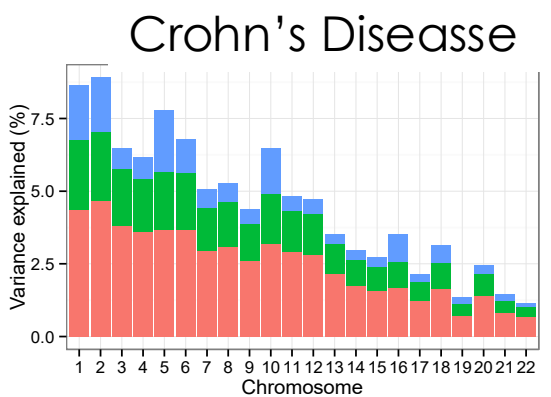
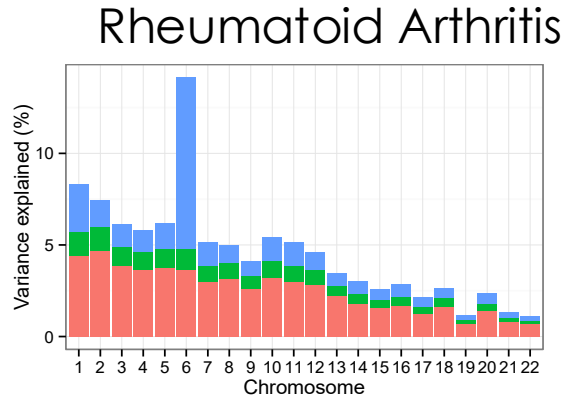
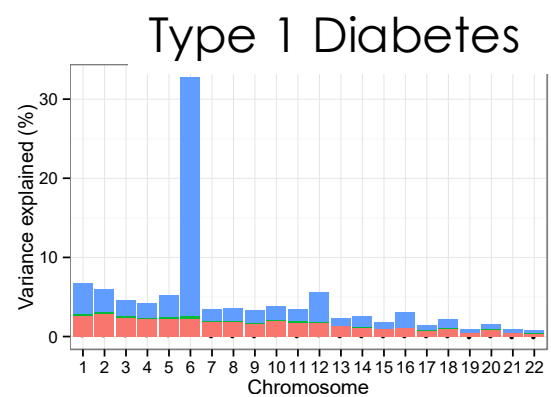
[Vassily Trubetskov](#), [Antonio F. Pardiñas](#), [Ting Qi](#), [Georgia Panagiotaropoulou](#), [Swapnil Awasthi](#), [Tim B. Bigdeli](#), [Julien Bryois](#), [Chia-Yen Chen](#), [Charlotte A. Dennison](#), [Lynsey S. Hall](#), [Max Lam](#), [Kyoko Watanabe](#), [Oleksandr Frej](#), [Tian Ge](#), [Janet C. Harwood](#), [Frank Koopmans](#), [Sigurdur Magnusson](#), [Alexander L. Richards](#), [Julia Sidorenko](#), [Yang Wu](#), [Jian Zeng](#), [Jakob Grove](#), [Minsoo Kim](#), [Zhiqiang Li](#), [Indonesia Schizophrenia Consortium](#), [PsychENCODE](#), [Psychosis Endophenotypes International Consortium](#), [The SynGO Consortium](#), [Schizophrenia Working Group of the Psychiatric Genomics Consortium](#)

[+ Show authors](#)

[Nature](#) **604**, 502–508 (2022) | [Cite this article](#)

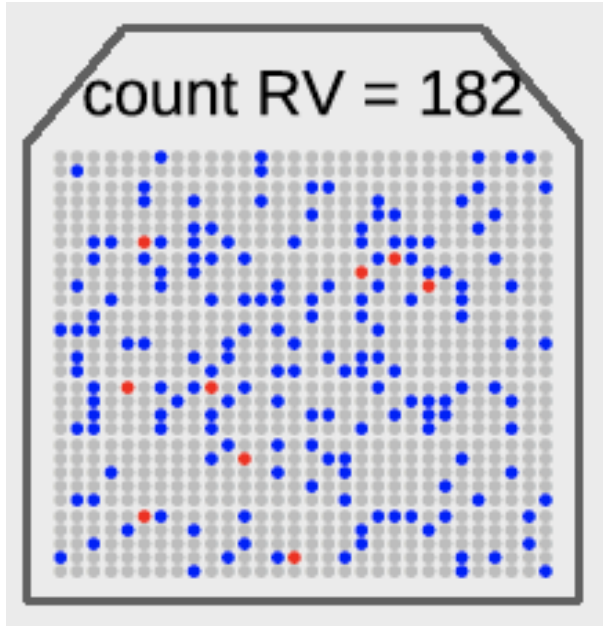
57k Accesses | **321** Citations | **463** Altmetric | [Metrics](#)





Many DNA variants contribute to genetic risk, and most have very small effects.

Polygenic disease for an individual



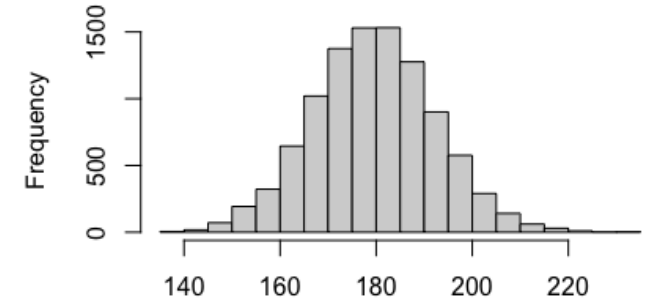
900 DNA polymorphic sites

RV = risk variant

Frequency of risk variant at each site: 0.1 (p)

Average person $900 \times 2 \times 0.1 = 180$ risk variant

Mean \pm 3SD: 142 to 218

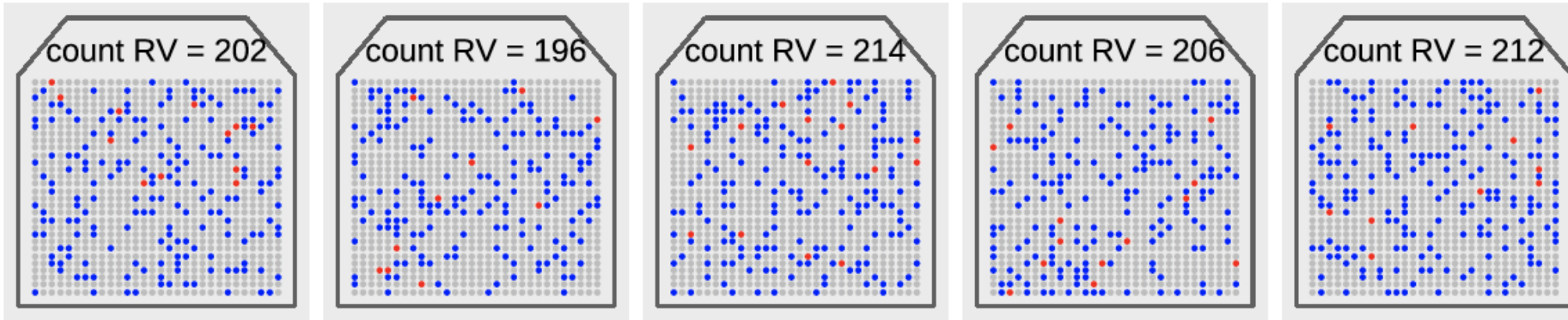


- 0 Grey: Homozygote no risk alleles (or equivalently 2 protective alleles)
- 1 Blue : Heterozygote one risk allele (and one non-risk/protective allele)
- 2 Red: Homozygote two risk alleles

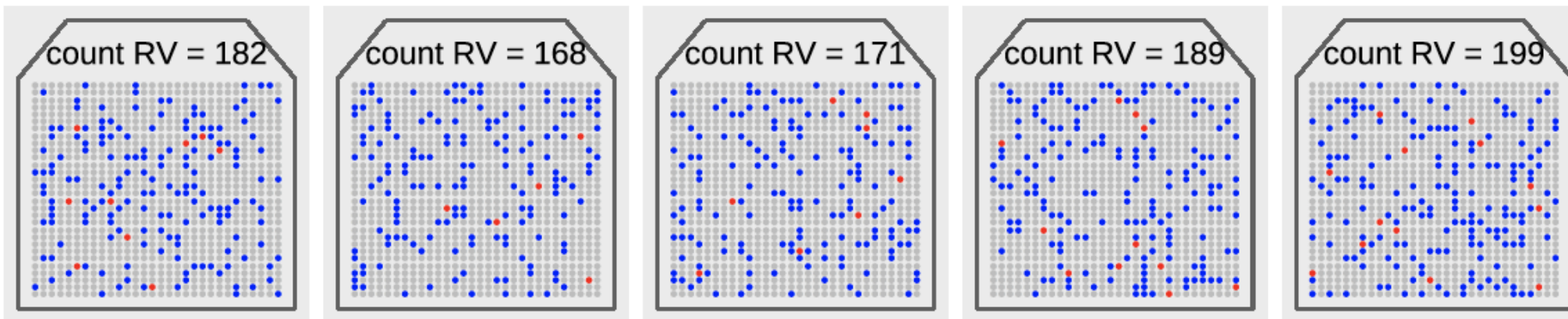
Toy
example

Polygenic disease for an individual

Affected over lifetime

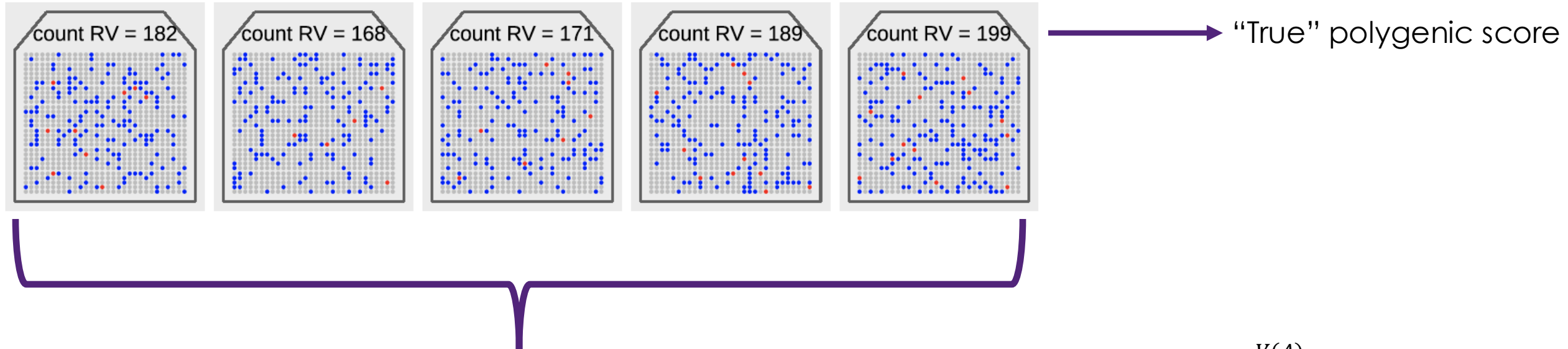


Not affected over lifetime



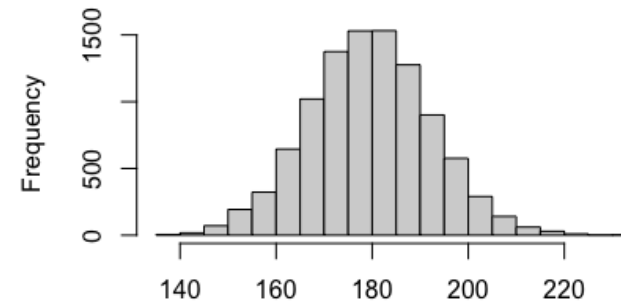
- We all carry risk variants for all diseases.
- Robustness
- Those affected carry a higher burden.
- Non-genetic factors contribute to risk too
- Each person carries a unique portfolio of risk alleles

Polygenic score

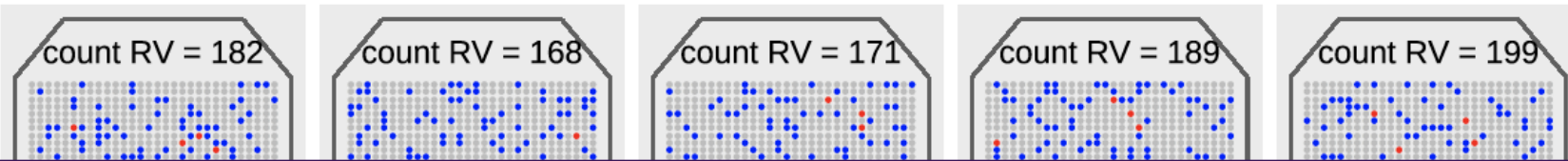


Genetic variance between people attributed to all genetic factors $V(A)$

$$h^2 = \frac{V(A)}{V(P)} \text{ heritability}$$



Polygenic score



→ “True” polygenic score

Not all variants captured
on genotyping arrays

Genetic variance between people attributed to all genetic factors $V(A)$

$$h^2 = \frac{V(A)}{V(P)} \text{ heritability}$$

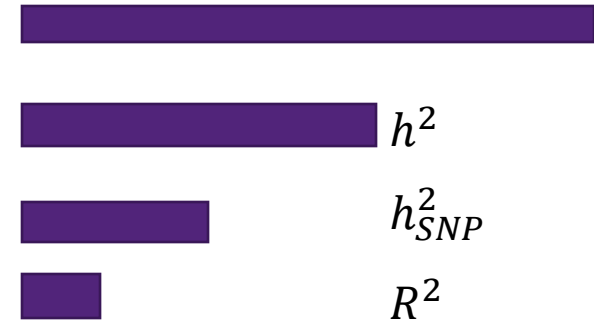
Genetic variance between people attributed to all genetic factors
associated with SNPs on genotyping arrays

$$h_{SNP}^2 = h_g^2 = \frac{V(A:SNP)}{V(P)}$$

SNP – based heritability

Limitations in prediction accuracy

- ❖ PGS have a **theoretical** upper limit dependent on the **heritability of the trait** (how much of the variance of trait values between people is attributed to genetic factors).
- ❖ PGS have a **technical** upper limit associated with the proportion of **variance tagged** by the DNA variants measured.
- ❖ PGS have a **practical** upper limit dependent on the **sample size of the discovery sample** used to estimate effect sizes of risk alleles, and the **quality** of the discovery sample.
- ❖ PGS can be pushed closer to the technical upper limit by the **statistical methodology** used to generate the optimal weighting given to the risk alleles, and new methods integrate new biological data.



Schizophrenia

Max:

25% Liability

AUC 0.84

Current:

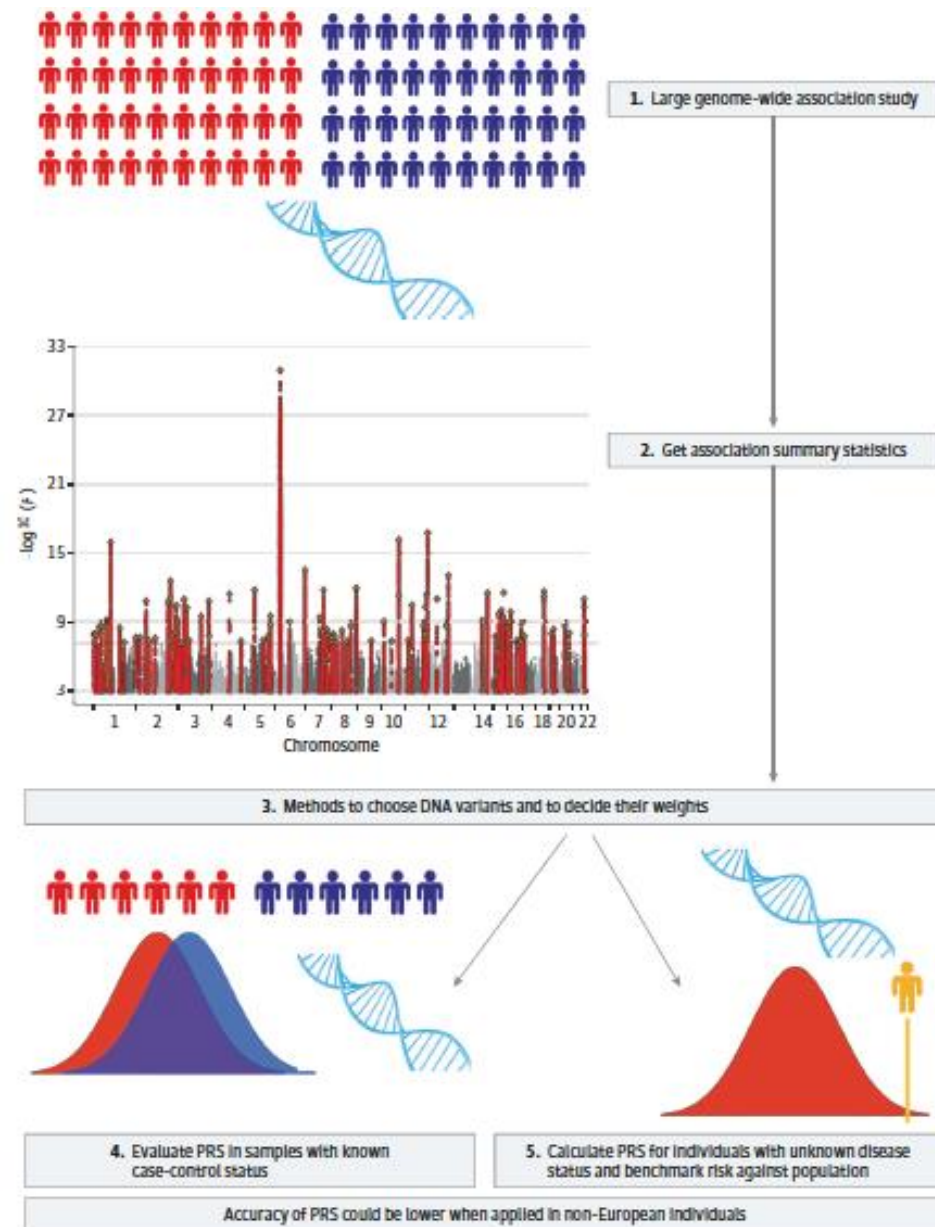
11% Liability

AUC 0.74

Polygenic scores cannot
be highly accurate
predictors of phenotypes

Basic method for PGS

Polygenic scores (PGS)



- A weighted count of risk alleles

$$\text{PGS} = \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \widehat{\beta}_3 x_{i3} + \dots = \sum_{j=1}^{n_{\text{SNP}}} \widehat{\beta}_j x_{ij}$$

0, 1 or 2
Risk alleles

Which SNPs?

What weights?

- Don't need to know causal variants for prediction!
- Prediction can be based on correlated variants.

4. Evaluate

$$Y = b \cdot \text{PGS} + e$$

$$R^2 = \text{var}(b \cdot \text{PGS}) / \text{Var}(Y)$$

AUC statistic:

Probability that a case ranks higher than a control

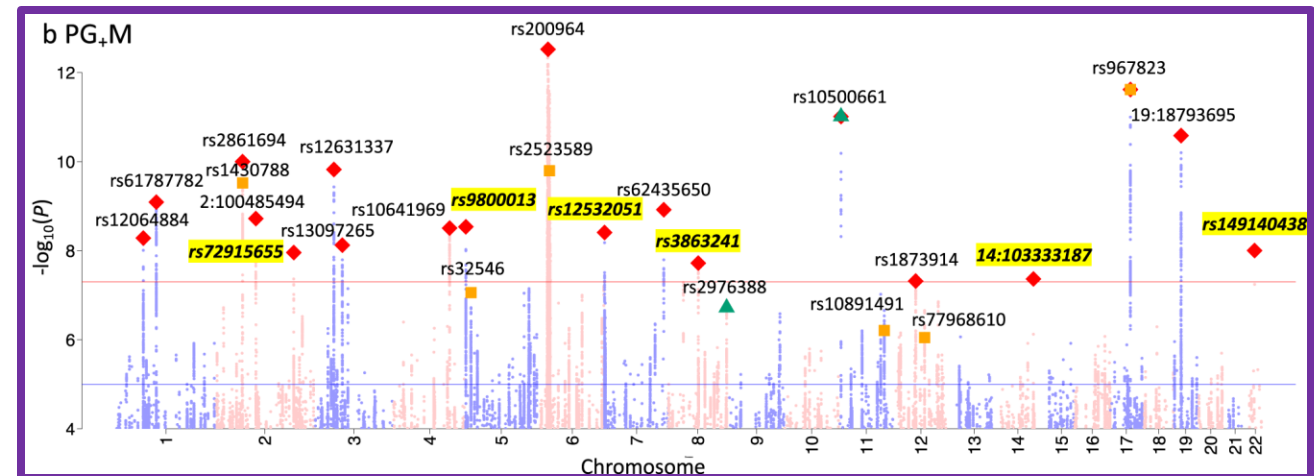
A weighted sum of the count of risk alleles

$$\text{PGS} = \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \widehat{\beta}_3 x_{i3} + \dots = \sum_{j=1}^{n_{\text{SNP}}} \widehat{\beta}_j x_{ij}$$

How many SNPs?
Which SNPs?
What weights?

Basic method:

Clumping & P-value thresholding
(C+PT):



A weighted sum of the count of risk alleles

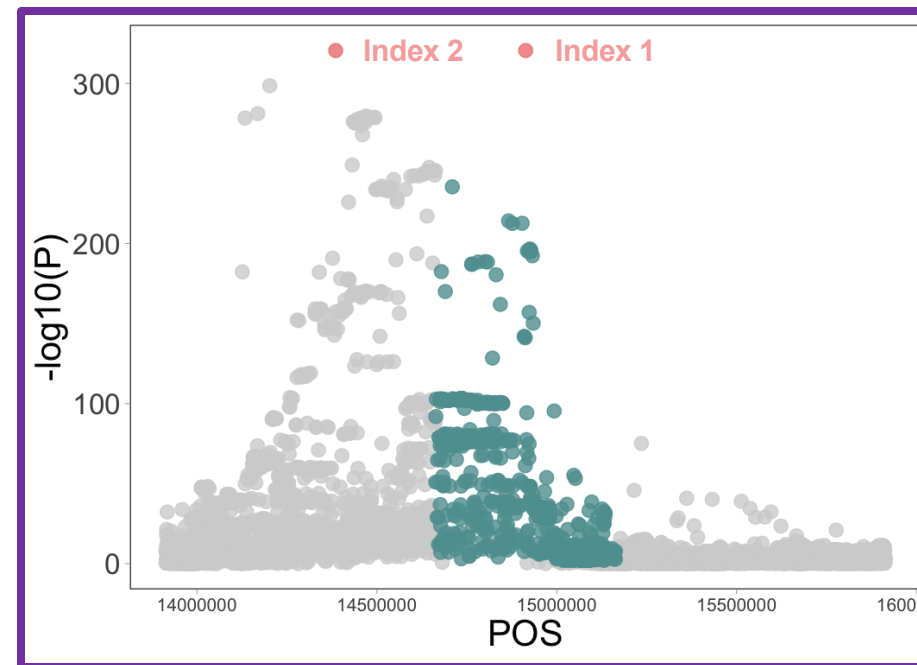
$$\text{PGS} = \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \widehat{\beta}_3 x_{i3} + \dots = \sum_{j=1}^{n_{\text{SNP}}} \widehat{\beta}_j x_{ij}$$

How many SNPs?
Which SNPs?
What weights?

Basic method:

Clumping & P-value thresholding (C+PT):

- Select most associated SNP in tower – LD-based clumping
- Select on a p-value threshold



A weighted sum of the count of risk alleles

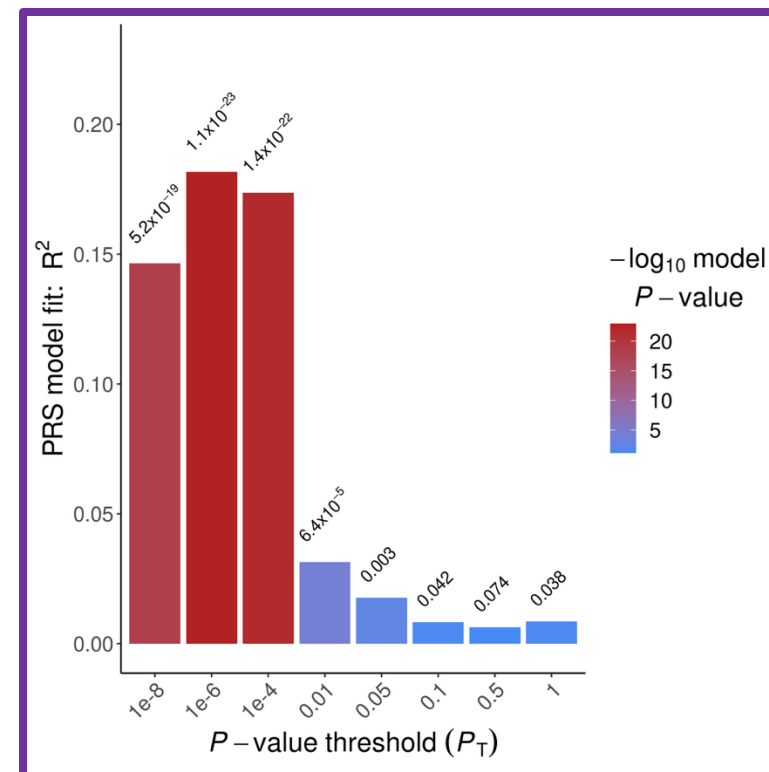
$$\text{PGS} = \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \widehat{\beta}_3 x_{i3} + \dots = \sum_{j=1}^{n_{\text{SNP}}} \widehat{\beta}_j x_{ij}$$

How many SNPs?
Which SNPs?
What weights?

Basic method:

Clumping & P-value thresholding (C+PT):

- Select most associated SNP in tower – LD-based clumping
- Select on a p-value threshold



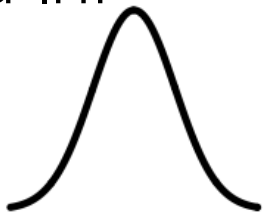
A weighted sum of the count of risk alleles

$$\text{PGS} = \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3} + \dots = \sum_{j=1}^{n_{\text{SNP}}} \hat{\beta}_j x_{ij}$$

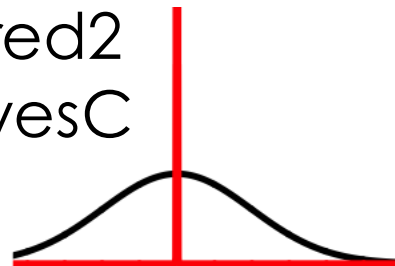
How many SNPs?
Which SNPs?
What weights?

**New methods model
genetic architecture**

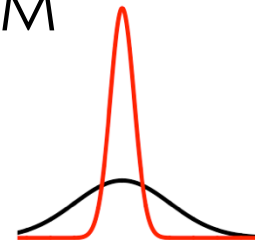
LDpred-Inf
SBLUP



LDPred2
SBayesC



BSLMM



SBayesR

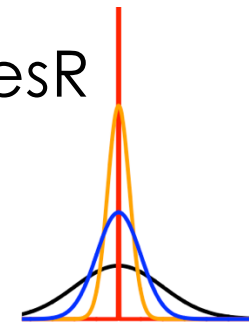


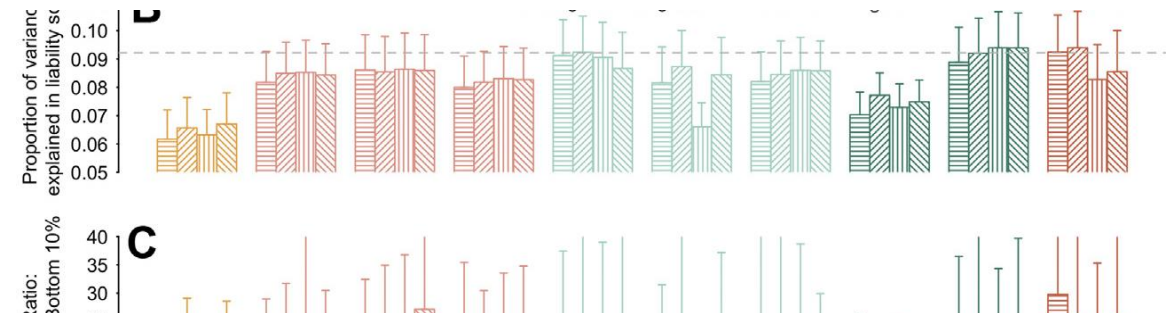
Table 1. Summary of Methods Used to Generate Polygenic Scores

Method	Distribution of SNP Effects (β)	Tuning Sample	Predefined Parameters	Parameters Estimated in Tuning Sample
PC+T	None	Yes	—	p -value threshold
SBLUP	$\beta \sim N\left(0, \frac{h_g^2}{m}\right)$ h_g^2 : SNP-based heritability, m : number of SNPs; $\lambda = m(1 - h_g^2)/h_g^2$	No	λ LD radius in kb	—
Ldpred2-Inf	Same as SBLUP	No	h_g^2 LD radius in cM or kb	—
Ldpred-funct	$\beta_j \sim N(0, c\sigma_j^2)$ $\sum_{j=1}^M 1_{\sigma_j^2 > 0} c\sigma_j^2 = h_g^2$, c is a normalizing constant, σ_j^2 is the expected per SNP heritability under the baseline-LD annotation model estimated by stratified LDSC from the discovery GWAS within Ldpred-funct software	No	h_g^2 LD radius in number of SNPs	—
Ldpred2	$\beta_j \sim \begin{cases} N\left(0, \frac{h_g^2}{\pi m}\right), & \text{with probability of } \pi \\ 0, & \text{with probability of } 1 - \pi \end{cases}$ When sparsity is “true,” the β_j for SNPs in the $(1 - \pi)$ partition are all set to zero	Yes	h_g^2 π software default values, LD radius in cM or kb	π , sparsity
Lassosum	$f(\beta) = \mathbf{y}^T \mathbf{y} + (1 - s)\beta^T \mathbf{X}^T \mathbf{X} \beta - 2\beta^T \mathbf{X}^T \mathbf{y} + s\beta^T \beta + 2\lambda \ \beta\ _1$ \mathbf{X} : $n \times m$ matrix of genotypes of LD reference sample, where n is sample size	Yes	LD blocks	λ , s
PRS-CS	$\beta_j \sim N\left(0, \frac{\sigma_j^2}{n}\psi_j\right)$ $\psi_j \sim G(a, \delta_j)$ $\delta_j \sim G(b, \phi)$, ϕ is a global scaling parameter	Yes	$a = 1, b = 0.5$ n LD blocks	ϕ
PRS-CS-auto	Same as PRS-CS, but estimates ϕ from the discovery GWAS	No	$a = 1, b = 0.5$ n LD blocks	—
SBayesR	$\beta_j \pi, \sigma_j^2 \sim \begin{cases} 0, & \text{with probability of } \pi_1 \\ N(0, \gamma_2 \sigma_j^2), & \text{with probability of } \pi_2 \\ \vdots \\ N(0, \gamma_c \sigma_j^2), & \text{with probability of } 1 - \sum_{c=1}^{C-1} \pi_c \end{cases}$ $\sigma_j^2 \sim \text{Inv-}\chi^2(d.f. = 4)$ $\pi_1 \sim \text{Dir}(1)$, estimated from discovery GWAS in SBayesR software γ_i are scaling parameters	No	LD radius in cM or kb $C = 4$ γ software default values	—
MegaPRS	Lasso: $\beta_j \sim DE(\lambda/\sigma_j)$ Ridge regression: $\beta_j \sim N(0, v\sigma_j^2)$ BOLT-LMM: $\beta_j \sim \begin{cases} N\left(0, \frac{(1-f_2)\sigma_j^2}{\pi}\right), & \text{with probability of } \pi \\ N\left(0, \frac{f_2\sigma_j^2}{1-\pi}\right), & \text{with probability of } 1 - \pi \end{cases}$ f_2 is the proportion of the total mixture variance in the second normal distribution BayesR: similar to SBayesR with $C = 4$, and π_i and γ_i estimated in the tuning sample σ_j^2 is the expected per SNP-heritability under BLD-LDAK model using SumHer	Yes	LD radius in cM or kb Parameters used in BLD-LDAK Grid search parameter values for each method	The tuning cohort is used to estimate the parameters that maximize prediction for each model, and from these the model that maximizes prediction is selected

Archival Report

A Comparison of Ten Polygenic Score Methods for Psychiatric Disorders Applied Across Multiple Cohorts

Guiyan Ni, Jian Zeng, Joana A. Revez, Ying Wang, Zhili Zheng, Tian Ge, Restuadi Restuadi, Jacqueline Kiewa, Dale R. Nyholt, Jonathan R.I. Coleman, Jordan W. Smoller, Schizophrenia Working Group of the Psychiatric Genomics Consortium, Major Depressive Disorder Working Group of the Psychiatric Genomics Consortium, Jian Yang, Peter M. Visscher, and Naomi R. Wray



- Random effects models > fixed effects models
- Mixture models > non-mixture (infinitesimal) models

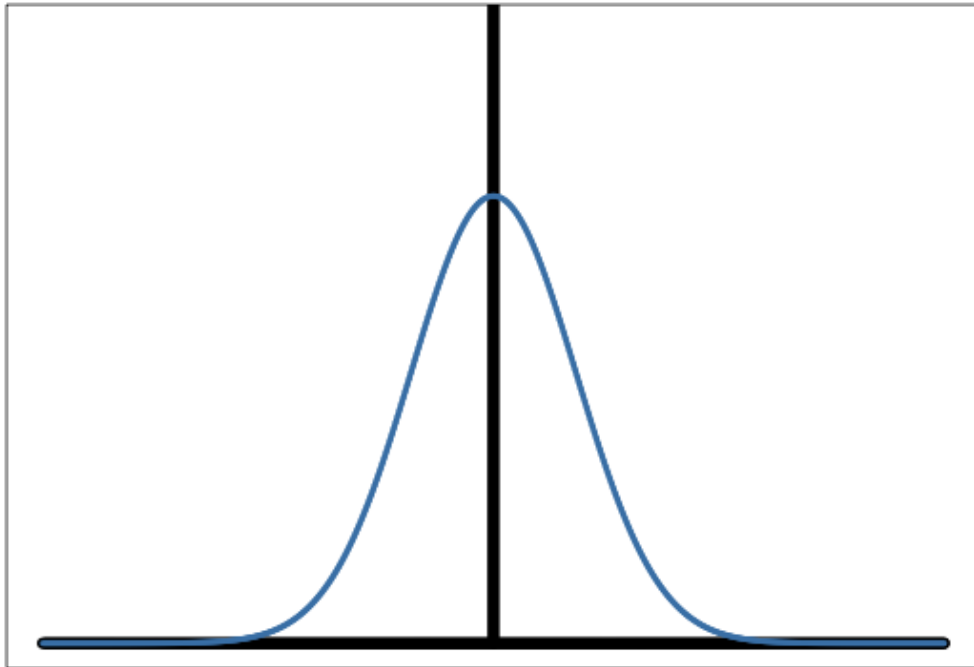
Bayesian methods for PGS

- Bayesian methods can estimate all parameters including SNP effects simultaneously
- Allow alternative assumptions regarding the distribution of SNP effects
- Different Bayesian methods mainly differ in
 - Assumption of the SNP effect distribution
 - Implementation

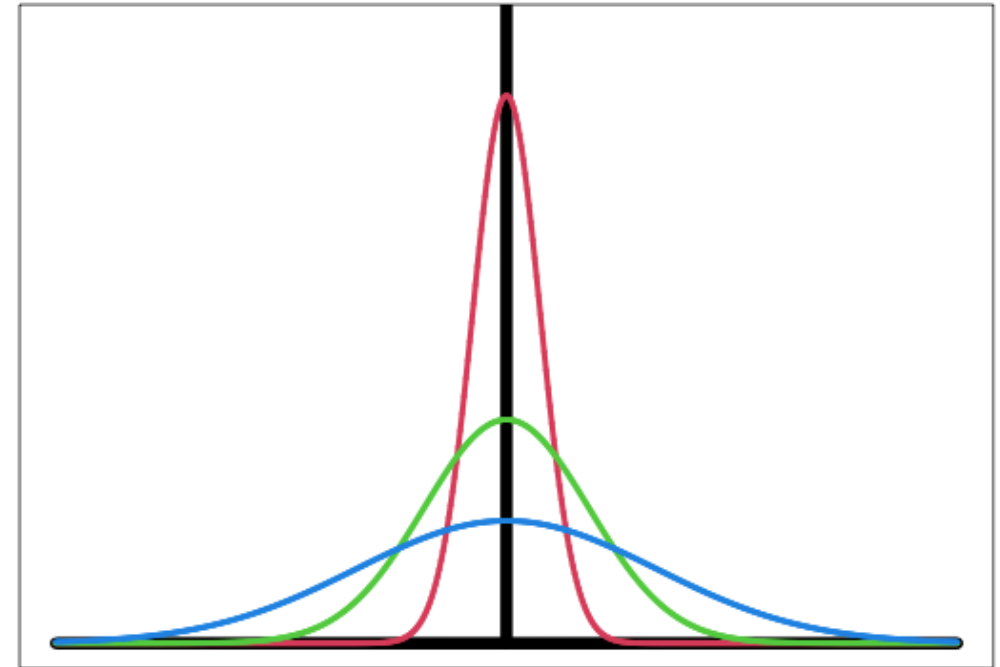
Assumptions for SNP effect distribution

Assumption	Distribution of SNP effects	Method
Infinite number of small effects	Normal distribution	BLUP/LDpred2-inf
Small number of moderate to large effects, many small effects	Student's t distribution, hierarchical distribution	BayesA/PRS-CS
Small number of small effects, many zero effects	Mixture, spike at zero, normal distribution	LDpred2/BayesC
Many zero effects, proportion of small effects, some moderate to large effects	Mixture, multiple normals	BayesR

BayesC



BayesR



How to incorporate this prior knowledge in the estimation of SNP effects?

Bayes theorem

$$P(x | y) \propto P(y | x)P(x)$$

Probability of
parameters x given
the data y (**posterior**)

Is proportional to

Probability of
data y given the
 x (**likelihood** of
data)

Prior
probability
of x

Consider an experiment where we measure height of 10 people
to estimate average height

We want to use prior knowledge from many previous studies that average height is 174cm with standard error 5cm

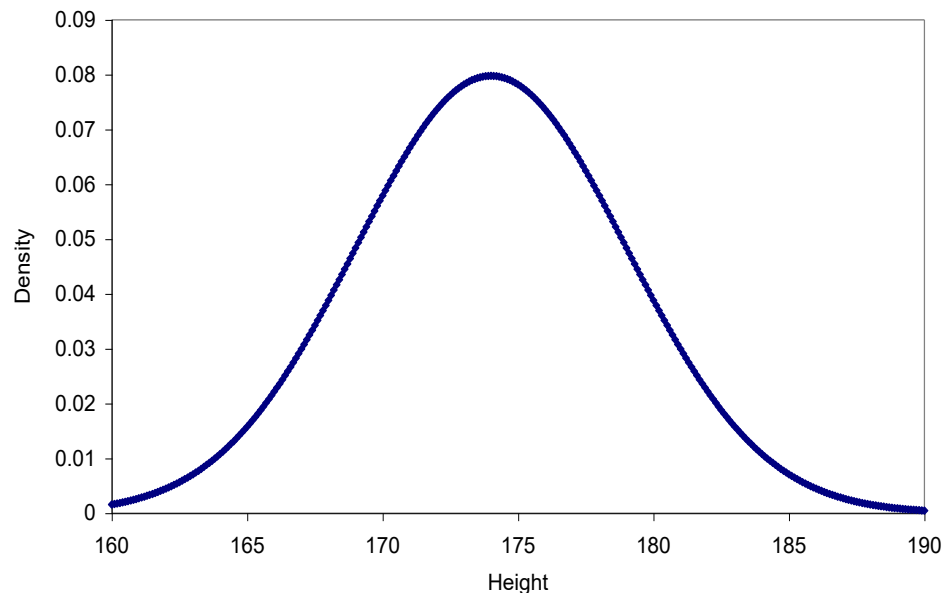
$$y = \text{average height} + e$$

Bayes theorem

$$P(x | y) \propto P(y | x)P(x)$$



Prior probability of x (average height)



Bayes theorem

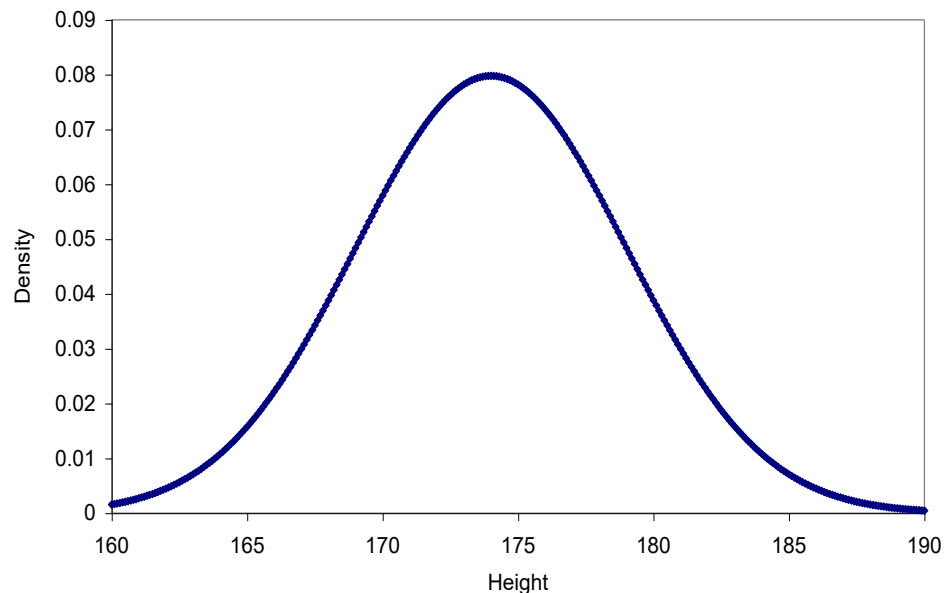
$$P(x | y) \propto P(y | x)P(x)$$

From the data.....

$$\bar{x} = 178$$

$$s.e = 5$$

Prior probability of x (average height)

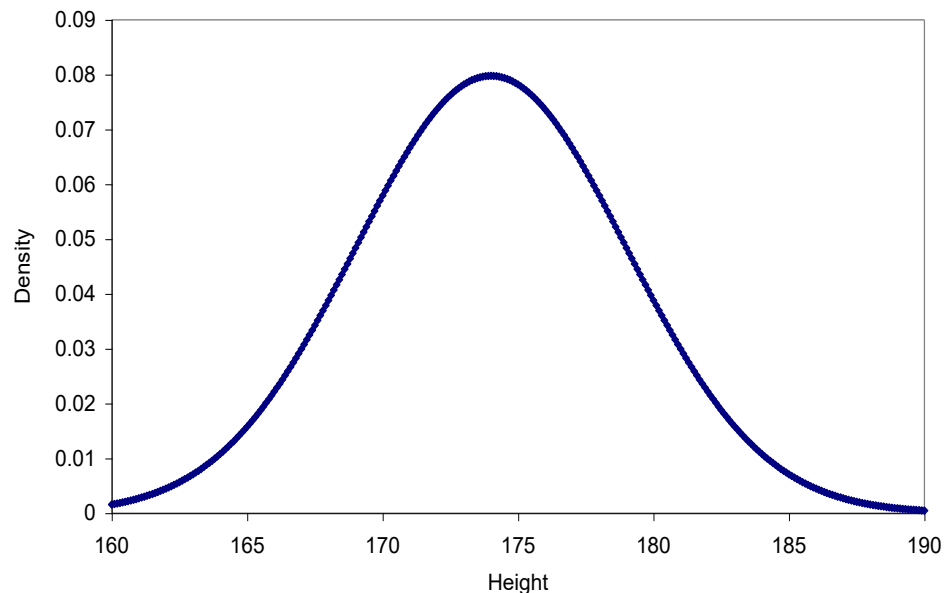
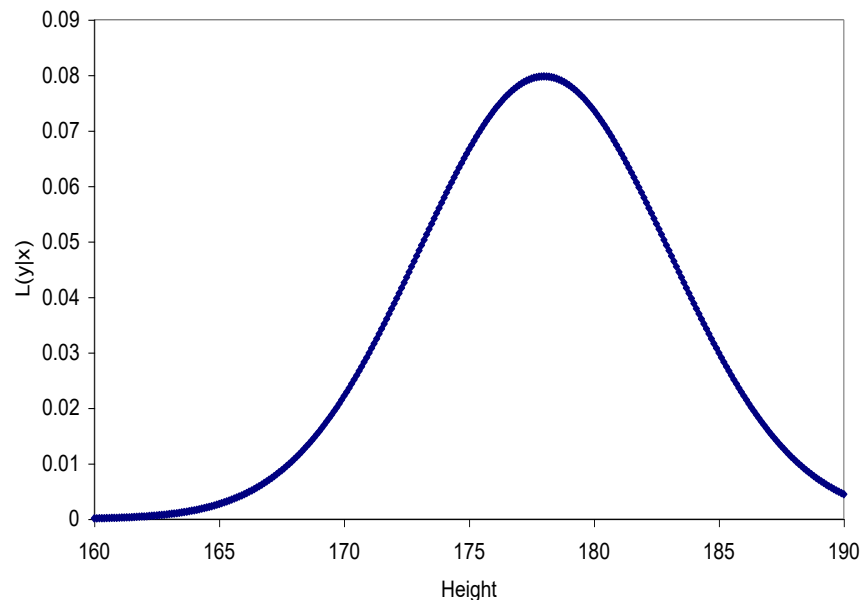


Bayes theorem

$$P(x | y) \propto P(y | x)P(x)$$

Likelihood of data (y) given
height x, most likely x = 178cm

Prior probability of x (average height)



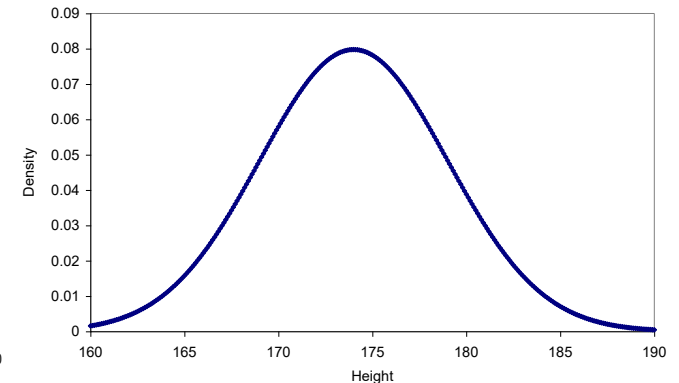
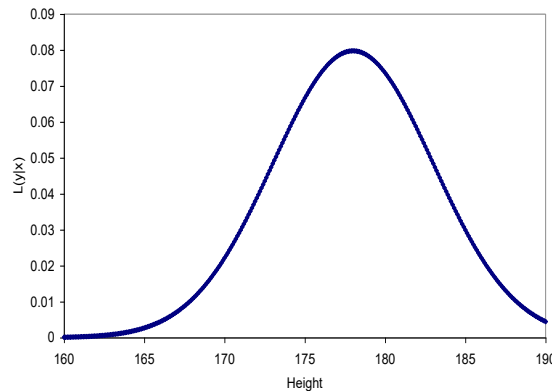
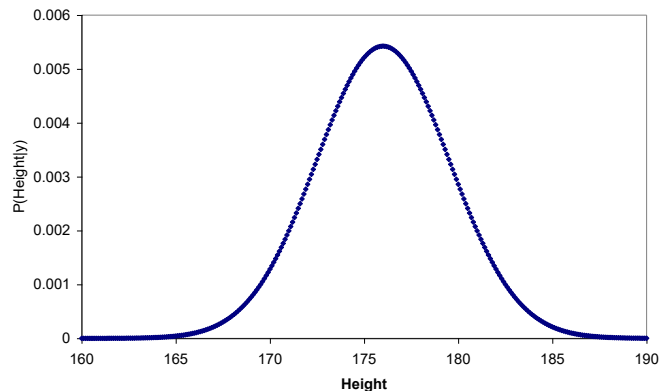
Bayes theorem

$$P(x | y) \propto P(y | x)P(x)$$

$P(x|y)$ mean = 176cm

$L(y|x)$

$P(x)$



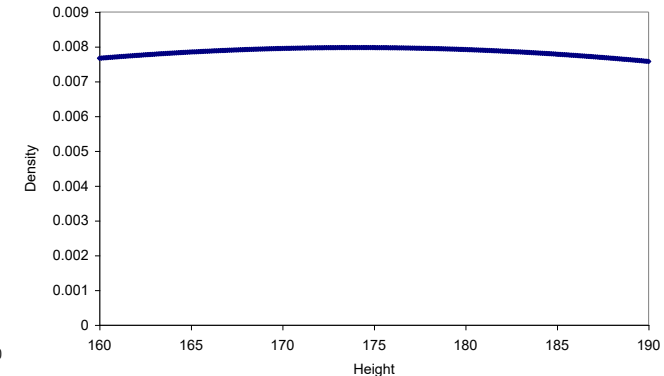
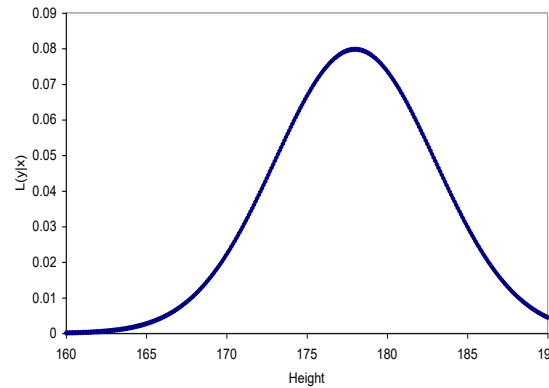
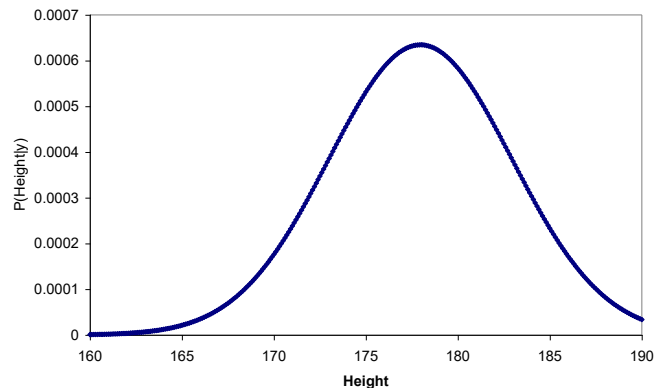
Less certainty about prior information? Use *less* informative (flat) prior

$$P(x | y) \propto P(y | x)P(x)$$

$P(x|y)$ mean = 178cm

$L(y|x)$

$P(x)$



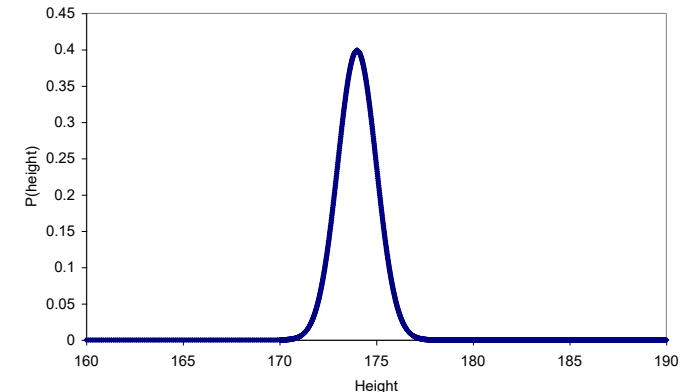
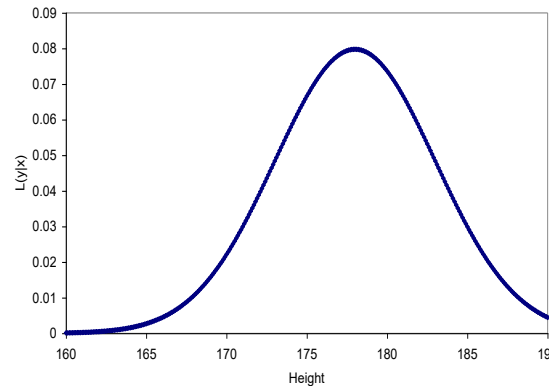
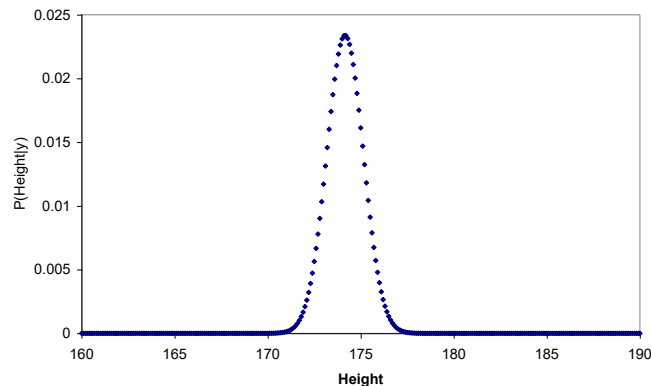
More certainty about prior information? Use *more* informative prior

$$P(x | y) \propto P(y | x)P(x)$$

$P(x|y)$ mean = 174.5cm

$L(y|x)$

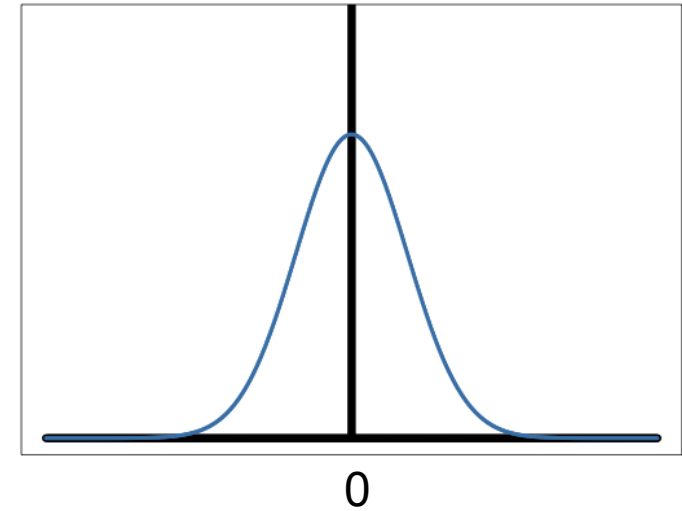
$P(x)$



Model

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

$$\beta_j \begin{cases} \sim N(0, \sigma_\beta^2) & \text{with probability } \pi \\ = 0 & \text{with probability } 1 - \pi \end{cases}$$



Posterior inference on SNP effects

$$P(\boldsymbol{\beta}|\mathbf{y}) \propto P(\mathbf{y}|\boldsymbol{\beta})P(\boldsymbol{\beta})$$

$$\propto (\sigma_e^2)^{-\frac{n}{2}} \exp\left\{-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma_e^2}\right\} \prod_{j=1}^m \left[(\sigma_\beta^2)^{-\frac{1}{2}} \exp\left\{-\frac{\beta_j^2}{2\sigma_\beta^2}\right\} \pi + \varphi_0(1 - \pi) \right]$$

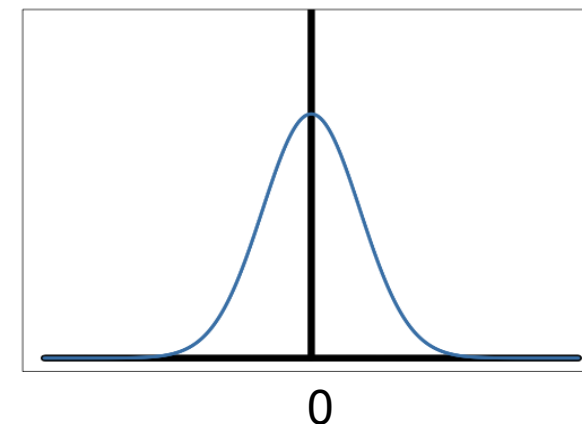
SNP effect estimates:

$$\hat{\boldsymbol{\beta}} = E(\boldsymbol{\beta}|\mathbf{y}) = \int \boldsymbol{\beta} P(\boldsymbol{\beta}|\mathbf{y}) d\boldsymbol{\beta}$$

$$= \int_{\beta_1} \dots \int_{\beta_m} (\sigma_e^2)^{-\frac{n}{2}} \exp\left\{-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma_e^2}\right\} \prod_{j=1}^m \left[(\sigma_\beta^2)^{-\frac{1}{2}} \exp\left\{-\frac{\beta_j^2}{2\sigma_\beta^2}\right\} \pi + \varphi_0(1 - \pi) \right] d\beta_1 \dots d\beta_m$$

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

$$\beta_j \begin{cases} \sim N(0, \sigma_\beta^2) & \text{with probability } \pi \\ = 0 & \text{with probability } 1 - \pi \end{cases}$$



Posterior inference on SNP effects

$$\hat{\boldsymbol{\beta}} = E(\boldsymbol{\beta}|\mathbf{y}) = \int_{\beta_1} \dots \int_{\beta_m} (\sigma_e^2)^{-\frac{n}{2}} \exp\left\{-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma_e^2}\right\} \prod_{j=1}^m \left[(\sigma_{\beta}^2)^{-\frac{1}{2}} \exp\left\{-\frac{\beta_j^2}{2\sigma_{\beta}^2}\right\} \pi + \varphi_0(1 - \pi) \right] d\beta_1 \dots d\beta_m$$

- Cannot solve directly \rightarrow no closed form solution
- Estimates of parameters depend on other parameters
- Use Markov chain Monte Carlo (MCMC) algorithm!

Markov chain

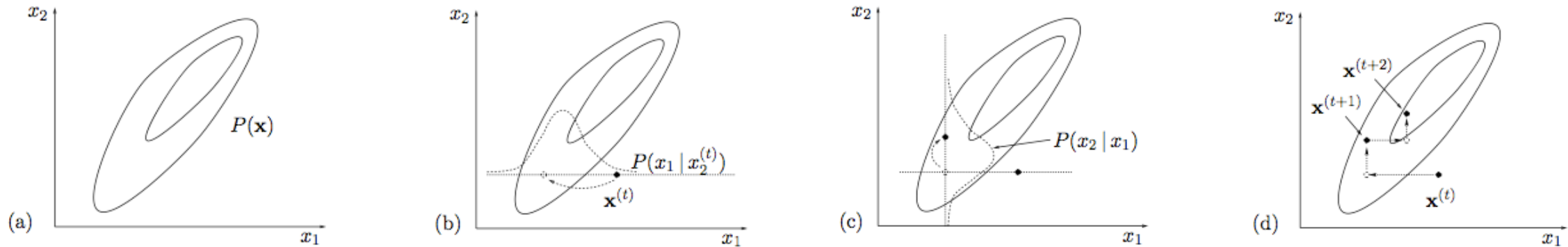
A sequence of samples where each sample depends only on the previous one (memoryless). This property allows the algorithm to gradually explore the distribution.

Monte Carlo

Using random sampling to perform numerical estimation, e.g., integrating over a probability distribution by averaging over samples.

Gibbs Sampling

A special case of MCMC to sample from posterior distribution of each parameter **conditional** on all other parameters.



The key is to derive $P(x_1 | x_2)$ and $P(x_2 | x_1)$

[Figure source](#)

To run Gibbs sampling, we need to derive the full conditional distribution for each parameter

- $P(\mu | \mathbf{y}, \boldsymbol{\beta}, \sigma_{\beta}^2, \pi, \sigma_e^2)$
- $P(\beta_j | \mathbf{y}, \boldsymbol{\beta}_{-j}, \sigma_{\beta}^2, \pi, \sigma_e^2)$
- $P(\sigma_{\beta}^2 | \mathbf{y}, \boldsymbol{\beta}, \pi, \sigma_e^2)$
- $P(\pi | \mathbf{y}, \boldsymbol{\beta}, \sigma_{\beta}^2, \sigma_e^2)$
- $P(\sigma_e^2 | \mathbf{y}, \boldsymbol{\beta}, \sigma_{\beta}^2, \pi)$

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

$$\beta_j \begin{cases} \sim N(0, \sigma_{\beta}^2) & \text{with probability } \pi \\ = 0 & \text{with probability } 1 - \pi \end{cases}$$

Gibbs sampling

- Set starting values for $(\mu, \boldsymbol{\delta}, \boldsymbol{\beta}, \sigma_{\beta}^2, \pi, \sigma_e^2)$
- Then (for many iterations)
 - For each SNP, sample δ_j, β_j conditional on other parameters
 - Sample $\mu, \sigma_{\beta}^2, \pi, \sigma_e^2$ with updated $\boldsymbol{\delta}, \boldsymbol{\beta}$

Samples reconstruct posterior distributions of parameters

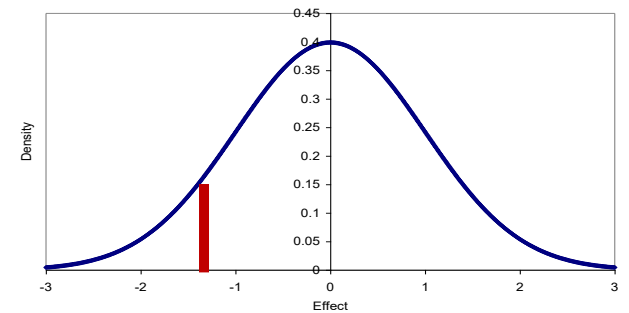
Gibbs Sampling

For example, for SNP effect β_j

- First sample if in zero or nonzero effect component of distribution (δ_j)
- Then if in nonzero component of the distribution, sample from

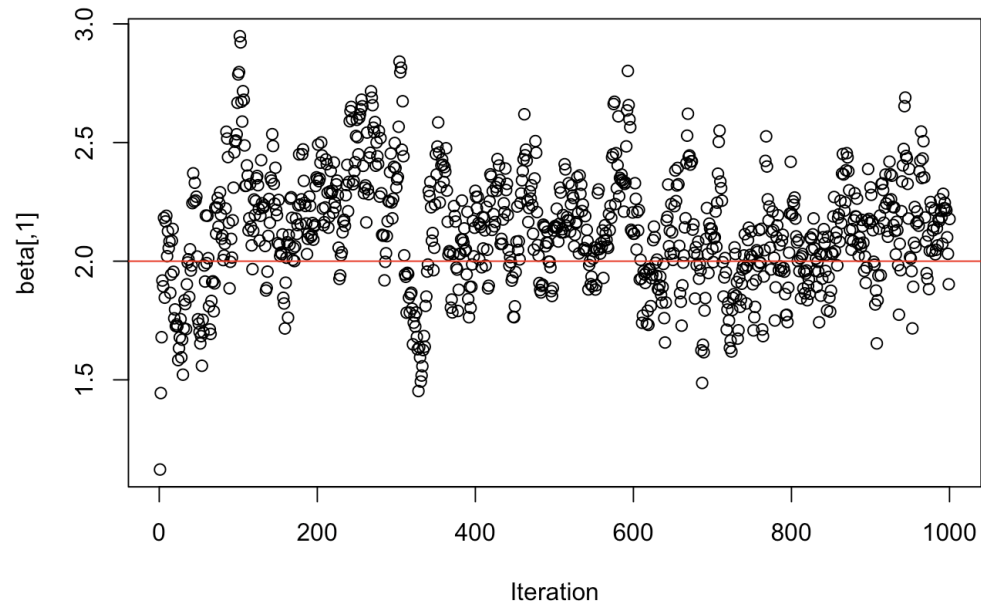
$$N\left(\frac{\mathbf{X}'_j(\mathbf{y} - \mathbf{1}_n\mu - \sum_{k \neq j} \mathbf{X}'_k \beta_k)}{\mathbf{X}'_j \mathbf{X}_j + \boxed{\sigma_e^2 / \sigma_\beta^2}}, \sigma_e^2 / (\mathbf{X}'_j \mathbf{X}_j + \sigma_e^2 / \sigma_\beta^2)\right)$$

Introduce shrinkage

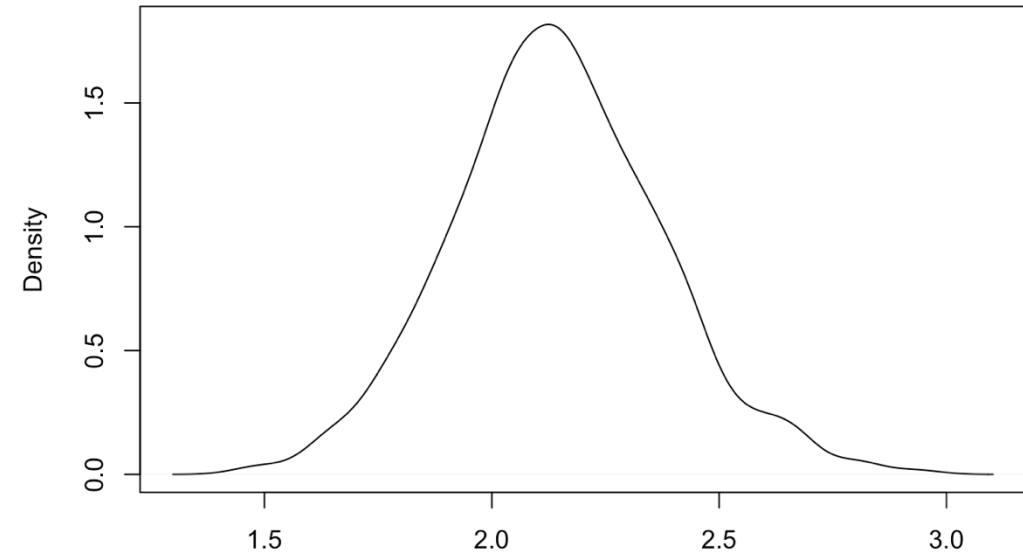


Gibbs sampling

Trace plot



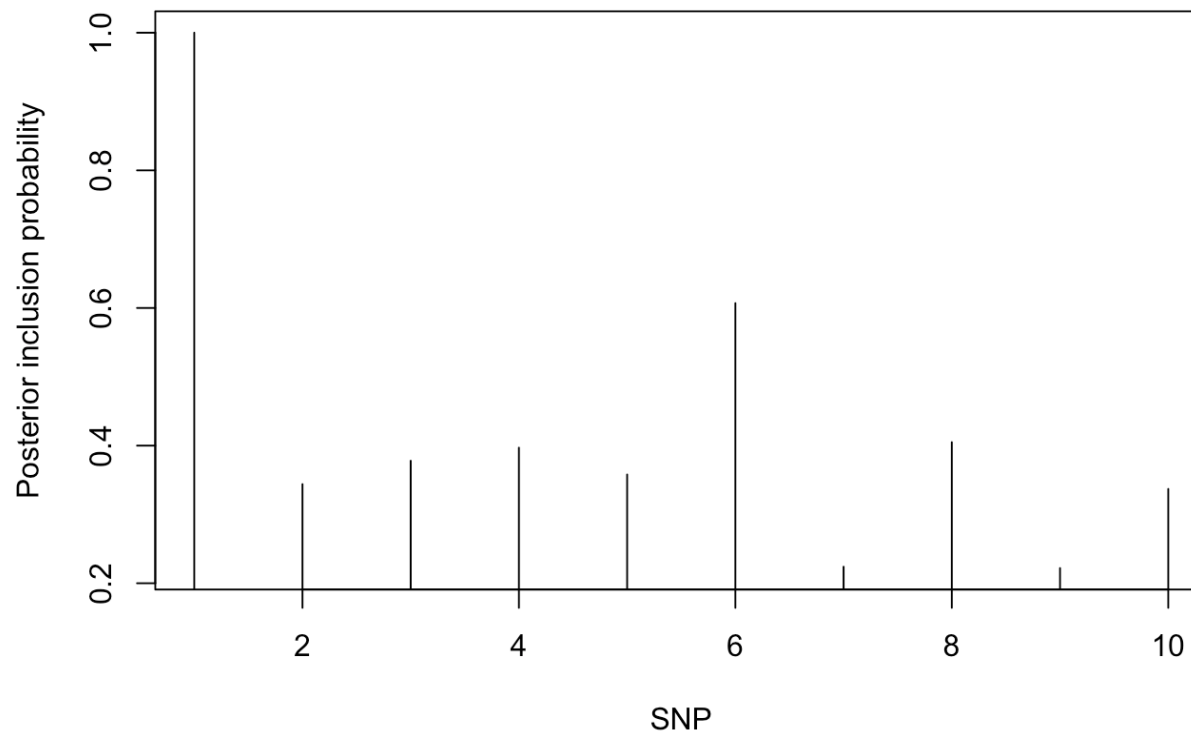
Posterior distribution



Posterior mean is used as the point estimate of the SNP effect

As a method of fine-mapping

Posterior inclusion probability (PIP):
probability that the SNP is included in the model with a nonzero effect.

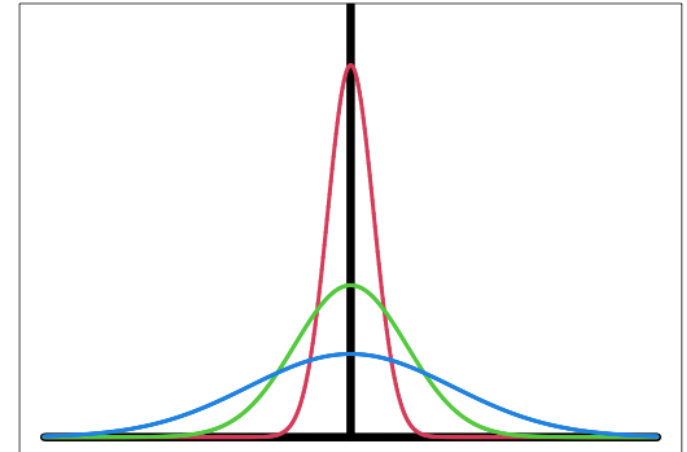


Model

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

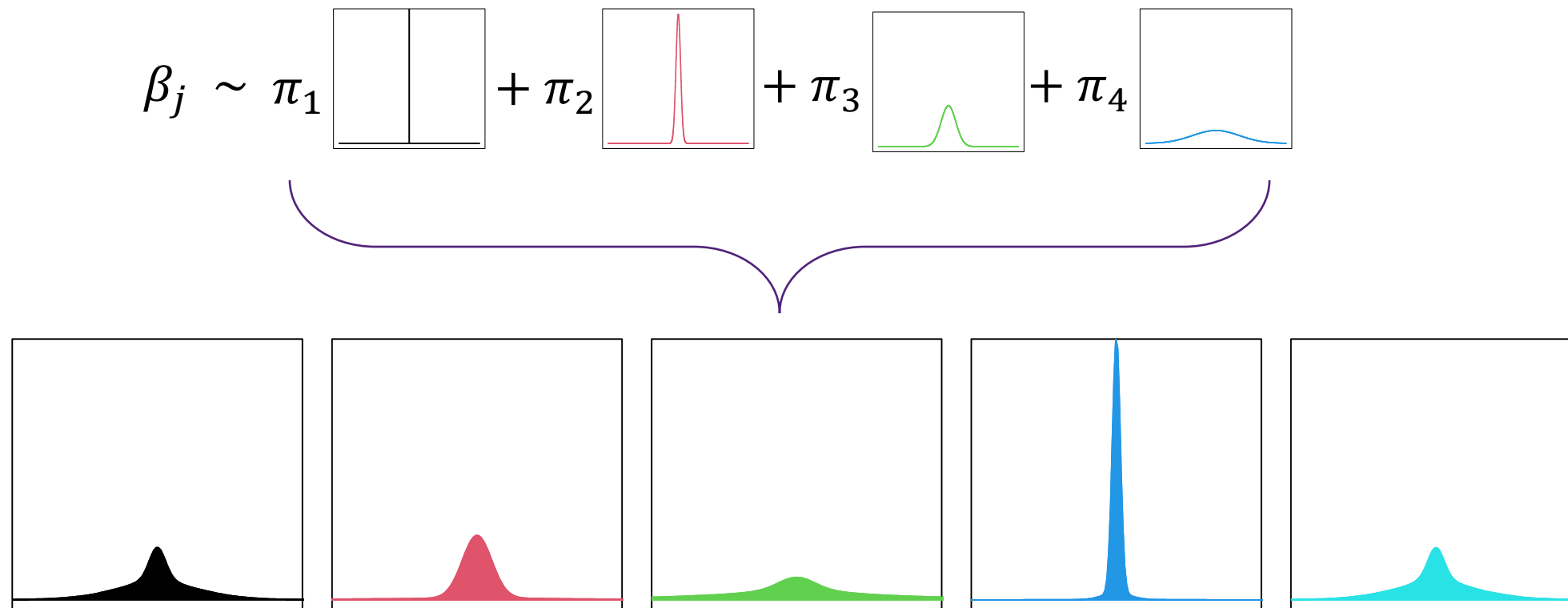
$$\beta_j | \pi, \sigma_\beta^2 = \begin{cases} 0 & \text{with probability } \pi_1, \\ \sim N(0, \gamma_2 \sigma_\beta^2) & \text{with probability } \pi_2, \\ \vdots & \\ \sim N(0, \gamma_C \sigma_\beta^2) & \text{with probability } 1 - \sum_{c=1}^{C-1} \pi_c, \end{cases}$$

$$\boldsymbol{\gamma} = (0, 0.01, 0.1, 1.0)'$$



BayesC is a special case of BayesR with two components

Why use multi-normal mixture?



Account for almost any distribution!

Practical Exercise Part 1: BayesR

<https://gctbhub.cloud.edu.au/data/teaching/>

Bayesian methods using GWAS summary statistics (sumstats)

- Best prediction methods take genetic values as random effect.
- These methods require individual genotypes and phenotypes.
- These data are often not publicly accessible.
- Computationally demanding with large # individuals/SNPs.
- Could be addressed by using GWAS summary statistics (**sumstats**).
- Methodology in human genetics has moved forward to use GWAS sumstats only.

Perspective

Workshop proceedings: GWAS summary statistics standards and sharing

2021

Jacqueline A.L. MacArthur,^{1,2,*} Annalisa Buniello,¹ Laura W. Harris,¹ James Hayhurst,¹ Aoife McMahon,¹ Elliot Sollis,¹ Maria Cerezo,¹ Peggy Hall,³ Elizabeth Lewis,¹ Patricia L. Whetzel,¹ Orli G. Bahcall,⁴ Inês Barroso,⁵ Robert J. Carroll,⁶ Michael Inouye,^{7,8,9} Teri A. Manolio,³ Stephen S. Rich,¹⁰ Lucia A. Hindorff,³ Ken Wiley,³ and Helen Parkinson^{1,*}

Table 1. Recommended standard reporting elements for GWAS SumStats

Data element	Column header	Mandatory/Optional
variant id	variant_id	One form of variant ID is mandatory, either rsID or chromosome, base pair location, and genome build ^a
chromosome	chromosome	
base pair location	base_pair_location	
p value	p_value	Mandatory
effect allele	effect_allele	Mandatory
other allele	other_allele	Mandatory
effect allele frequency	effect_allele_frequency	Mandatory
effect (odds ratio or beta)	odds_ratio or beta	Mandatory
standard error	standard_error	Mandatory
upper confidence interval	ci_upper	Optional
lower confidence interval	ci_lower	Optional

Genome-wide association studies

Emil Uffelmann¹, Qin Qin Huang², Nchangwi Syntia Munung³, Jantina de Vries³, Yukinori Okada^{4,5}, Alicia R. Martin^{6,7,8}, Hilary C. Martin², Tuuli Lappalainen^{9,10,12} and Danielle Posthuma^{1,11} ✉

Table 3 | **Databases of GWAS summary statistics**

Database	Content
GWAS Catalog ¹¹⁰	GWAS summary statistics and GWAS lead SNPs reported in GWAS papers
GeneAtlas ⁸	UK Biobank GWAS summary statistics
Pan UKBB	UK Biobank GWAS summary statistics
GWAS Atlas ²⁷³	Collection of publicly available GWAS summary statistics with follow-up in silico analysis
FinnGen results	GWAS summary statistics released from FinnGen, a project that collected biological samples from many sources in Finland
dbGAP	Public depository of National Institutes of Health-funded genomics data including GWAS summary statistics
OpenGWAS database	GWAS summary data sets
Pheweb.jp	GWAS summary statistics of Biobank Japan and cross-population meta-analyses

For a comprehensive list of genetic data resources, see REF.¹³. GWAS, genome-wide association studies; SNP, single-nucleotide polymorphism.

What are the minimum data required?

Given the standard GWAS with genotypes being allelic counts (0/1/2), the minimum data required for PGS prediction include:

- SNP marginal effect estimates
 - Standard errors
 - GWAS sample size
- } GWAS sumstats
-
- LD correlations among SNPs → LD matrix

SNP marginal effect estimates

GWAS estimates effect of each SNP one at a time from single SNP regression, so the estimate is marginal to (unconditional on) other SNPs.

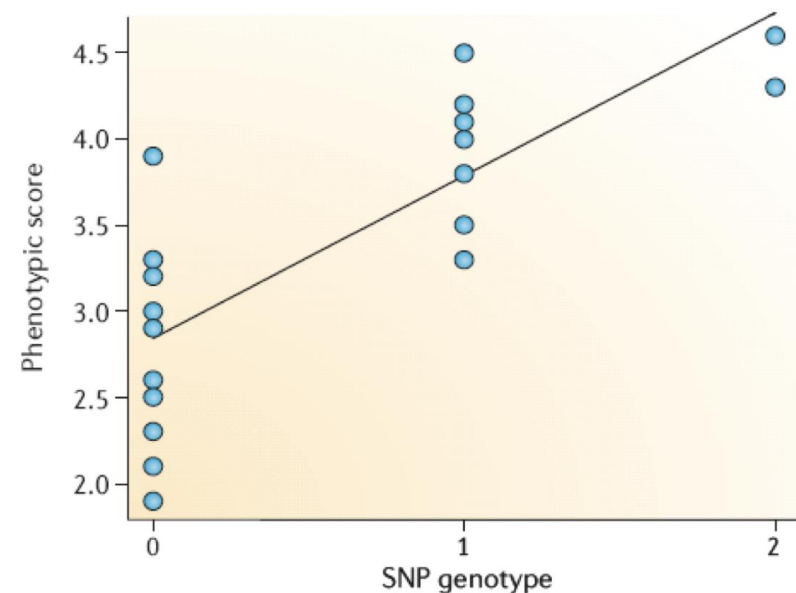
$$b_j = (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j \mathbf{y}$$

Assuming \mathbf{X} has been standardised with column mean zero and variance one, then

$$\mathbf{X}'_j \mathbf{X}_j = n \text{Var}(\mathbf{X}_j) = n$$

And

$$b_j = \frac{1}{n} \mathbf{X}'_j \mathbf{y}$$



Note that it has the inner product of the SNP genotypes and the phenotypes.

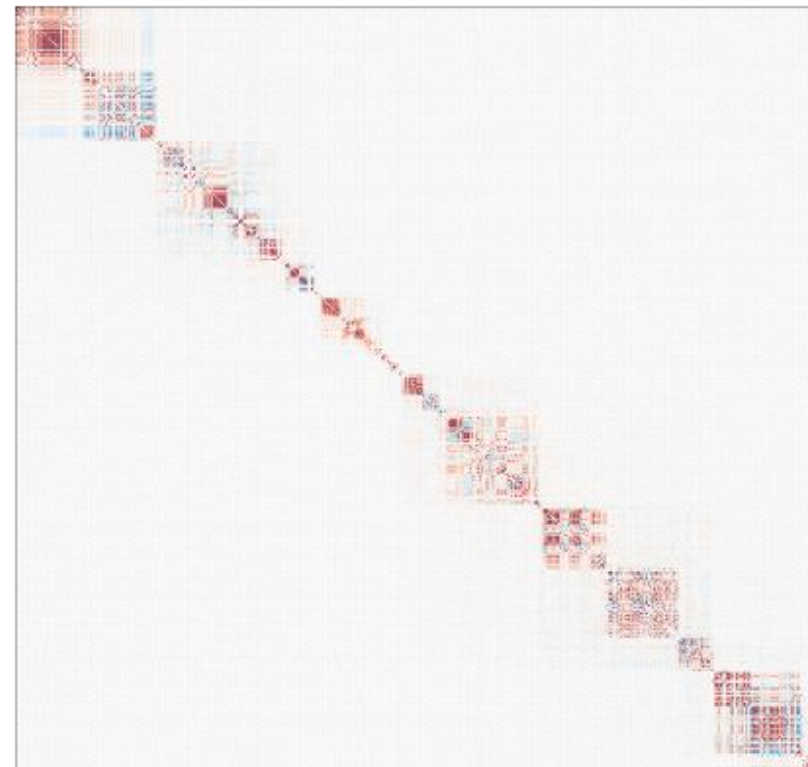
Linkage disequilibrium (LD) correlations

Usually obtained from a reference population

LD correlation matrix

$$\mathbf{R} = \frac{1}{n} \mathbf{X}'\mathbf{X}$$

assuming \mathbf{X} is standardised
with mean zero and
variance one

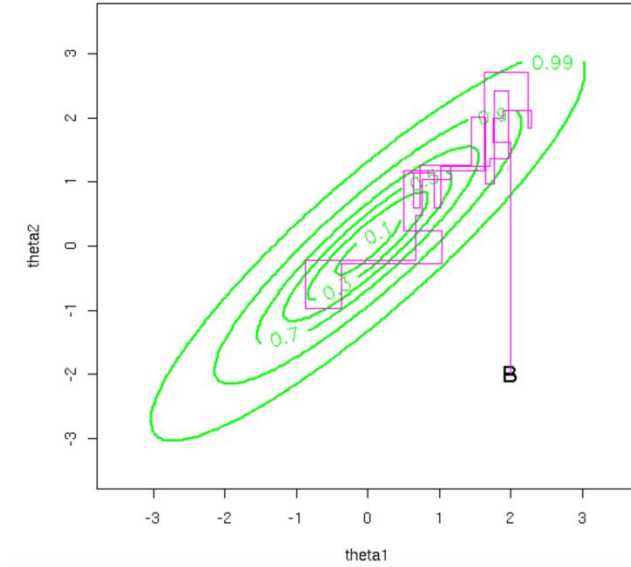


Gibbs sampling

Full conditional distribution for β_j , if in a nonzero dist'n,

$$f(\beta_j \mid \mathbf{b}, else) = N\left(\frac{r_j}{C_j}, \frac{\sigma_e^2}{C_j}\right)$$

where



Individual-level data

$$r_j = \mathbf{X}'_j \left(\mathbf{y} - \sum_{k \neq j} \mathbf{X}_k \beta_k \right)$$

$$C_j = \mathbf{X}'_j \mathbf{X}_j + \frac{\sigma_e^2}{\gamma_j \sigma_\beta^2}$$

Summary-level data

$$r_j = nb_j - \sum_{k \neq j} nR_{jk} \beta_k$$

$$C_j = n + \frac{\sigma_e^2}{\gamma_j \sigma_\beta^2}$$

Compare BayesR and SBayesR algorithms

All $\mathbf{X}'\mathbf{y}$ and $\mathbf{X}'\mathbf{X}$ can be replaced by $n\mathbf{b}$ and $n\mathbf{R}$

Algorithm 1 – Individual level data algorithm

```

Initialise parameters and read genotypes and phenotypes in PLINK binary format
Initialise  $\mathbf{y}^* = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$ 
for i := 1 to number of iterations do
  for j := 1 to p do
    Calculate  $r_j^* = \mathbf{x}_j' \mathbf{y}^*$ 
    Calculate  $r_j = r_j^* + \mathbf{x}_j' \mathbf{x}_j \beta_j^{(i-1)}$ 
    Calculate  $\sigma_c^2 = \sigma_\beta^2 \gamma_{\delta_j=c}$  for each of C classes (e.g., BayesR C=4 and  $\gamma = (0, 0.0001, 0.001, 0.01)$ )
    Calculate the left hand side  $l_{jc} = \mathbf{x}_j' \mathbf{x}_j + \frac{\sigma_c^2}{\sigma_\epsilon^2}$  for each of the C classes
    Calculate the log densities of given  $\delta_j = c$  using  $\log(\mathcal{L}_c) = -\frac{1}{2} \left[ \log \left( \frac{\sigma_c^2 l_{jc}}{\sigma_\epsilon^2} \right) - \frac{r_j^2}{\sigma_c^2 l_{jc}} \right] + \log(\pi_c)$ , where  $\pi_c$  is the current
    Calculate the full conditional posterior probability for  $\delta_j = c$  for C classes with  $\mathbb{P}(\delta_j = c | \boldsymbol{\theta}, \mathbf{y}) = \frac{1}{\sum_{l=1}^C \exp[\log(\mathcal{L}_l) - \log(\mathcal{L}_c)]}$ 
    Using full conditional posterior probabilities sample class membership for  $\beta_j^{(i)}$  using categorical random variable sampler
    Given class sample SNP effect  $\beta_j^{(i)}$  from  $N \left( \frac{r_j}{l_{jc}}, \frac{\sigma_c^2}{l_{jc}} \right)$ 
    Given SNP effect adjust corrected phenotype side  $(\mathbf{y}^*)^{(i)} = (\mathbf{y}^*)^{(i-1)} - \mathbf{x}_j \left( \beta_j^{(i)} - \beta_j^{(i-1)} \right)$ 
  od
od

Sample update from full conditional for  $\sigma_\beta^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\beta = v_\beta + q$  and  $\tilde{S}_\beta^2 = \frac{v_\beta S_\beta^2 + \sum_{j=1}^q \frac{\beta_j^2}{\gamma_c}}{v_\beta + q}$ ,
where  $q$  is the number of non-zero variants
Sample update from full conditional for  $\sigma_\epsilon^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\epsilon = n + v_\epsilon$ 
and scale parameter  $\tilde{S}_\epsilon^2 = \frac{SSE + v_\epsilon S_\epsilon^2}{n + v_\epsilon}$  and  $SSE = \mathbf{y}^*{}' \mathbf{y}^*$ 
Sample update from full conditional for  $\boldsymbol{\pi}$ , which is Dirichlet(C,  $\mathbf{c} + \boldsymbol{\alpha}$ ), where  $\mathbf{c}$  is a vector of length C and contains the counts
of the number of variants in each variance class and  $\boldsymbol{\alpha} = (1, \dots, 1)$ 
Calculate genetic variance for  $h_{SNP}^2$  calculation using  $\sigma_\delta^2 = \text{Var}(\mathbf{X}\boldsymbol{\beta})$ 
Calculate  $h_{SNP}^2 = \frac{\sigma_\delta^2}{\sigma_\delta^2 + \sigma_\epsilon^2}$ 
od

```

Algorithm 2 Summary data algorithm

```

Initialise parameters and read summary statistics
Reconstruct  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$  from summary statistics and LD reference panel
Calculate  $\mathbf{r}^* = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta}$ 
for i := 1 to number of iterations do
  for j := 1 to p do
    Calculate  $\mathbf{r}_j = \mathbf{r}_j^* + \mathbf{x}_j' \mathbf{x}_j \beta_j$ 
    Calculate  $\sigma_c^2 = \sigma_\alpha^2 \gamma_{\delta_j=c}$  for each of C classes (e.g., SBayesR C=4 and  $\gamma = (0, 0.01, 0.1, 1)'$ )
    Calculate the left hand side  $l_{jc} = \mathbf{x}_j' \mathbf{x}_j + \frac{\sigma_c^2}{\sigma_\epsilon^2}$  for each of the C classes
    Calculate the log densities of given  $\delta_j = c$  using  $\log(\mathcal{L}_c) = -\frac{1}{2} \left[ \log \left( \frac{\sigma_c^2 l_{jc}}{\sigma_\epsilon^2} \right) - \frac{r_j^2}{\sigma_c^2 l_{jc}} \right] + \log(\pi_c)$ , where  $\pi_c$  is the current
    Calculate the full conditional posterior probability for  $\delta_j = c$  for C classes with  $\mathbb{P}(\delta_j = c | \boldsymbol{\theta}, \mathbf{y}) = \frac{1}{\sum_{l=1}^C \exp[\log(\mathcal{L}_l) - \log(\mathcal{L}_c)]}$ 
    Using full conditional posterior probabilities sample class membership for  $\beta_j^{(i)}$  using categorical random variable sampler
    Given class sample SNP effect  $\beta_j^{(i)}$  from  $N \left( \frac{r_j}{l_{jc}}, \frac{\sigma_c^2}{l_{jc}} \right)$ 
    Given SNP effect adjust corrected right hand side  $(\mathbf{r}^*)^{(i+1)} = (\mathbf{r}^*)^{(i)} - \mathbf{X}'\mathbf{x}_j \left( \beta_j^{(i+1)} - \beta_j^{(i)} \right)$ .  $\mathbf{X}'\mathbf{x}_j$  is the  $j$ th column of  $\mathbf{X}'\mathbf{X}$ .
  od
od

Sample update from full conditional for  $\sigma_\alpha^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\alpha = v_0 + q$  and  $\tilde{\tau}_\alpha^2 = \frac{v_0 \tau_0^2 + \sum_{j=1}^q \frac{\beta_j^2}{\gamma_{\delta_j}}}{v_0 + q}$ ,
where  $q$  is the number of non-zero variants
Sample update from full conditional for  $\sigma_\epsilon^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\epsilon = n + v_\epsilon$ 
and scale parameter  $\tilde{\tau}_\epsilon^2 = \frac{SSE + v_\epsilon \tau_\epsilon^2}{n + v_\epsilon}$  and  $SSE = \mathbf{y}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{r}^* - \boldsymbol{\beta}'\mathbf{X}'\mathbf{y}$ 
Sample update from full conditional for  $\boldsymbol{\pi}$ , which is Dirichlet(C,  $\mathbf{c} + \boldsymbol{\alpha}$ ), where  $\mathbf{c}$  is a vector of length C and contains the counts
of the number of variants in each variance class.
Calculate genetic variance for  $h_{SNP}^2$  calculation using  $\sigma_\delta^2 = \hat{\boldsymbol{\beta}}' \mathbf{X}'\mathbf{y} - \hat{\boldsymbol{\beta}}' \mathbf{r}^*$ 
Calculate  $h_{SNP}^2 = \frac{\sigma_\delta^2}{\sigma_\delta^2 + \sigma_\epsilon^2}$ 
od

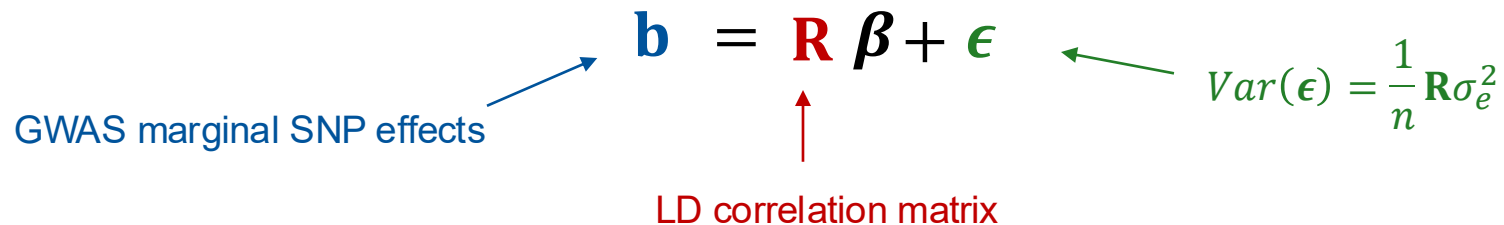
```

Consider an individual-data model with a standardised genotype matrix \mathbf{X} :

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

Multiply both sides by $\frac{1}{n}\mathbf{X}'$ gives

$$\frac{1}{n}\mathbf{X}'\mathbf{y} = \frac{1}{n}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \frac{1}{n}\mathbf{X}'\mathbf{e}$$


$$\mathbf{b} = \mathbf{R}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

GWAS marginal SNP effects

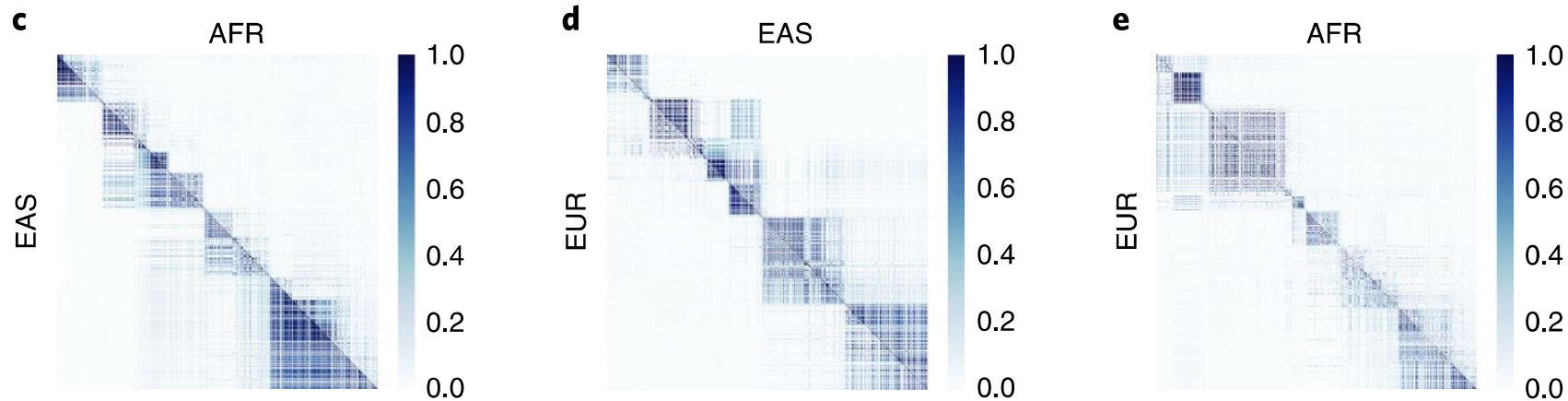
LD correlation matrix

$Var(\boldsymbol{\epsilon}) = \frac{1}{n}\mathbf{R}\sigma_e^2$

- In principle, SBayes and Bayes are equivalent methods when **same data** are used ($\mathbf{X}'\mathbf{y}$ and $\mathbf{X}'\mathbf{X}$ are sufficient statistics).
- However, when LD is estimated from a reference sample, SBayes is only an approximation to Bayes.
- Whether the difference is negligible depends on the heterogeneity in LD between the GWAS and LD ref samples.

LD reference population matches with GWAS population in genetics

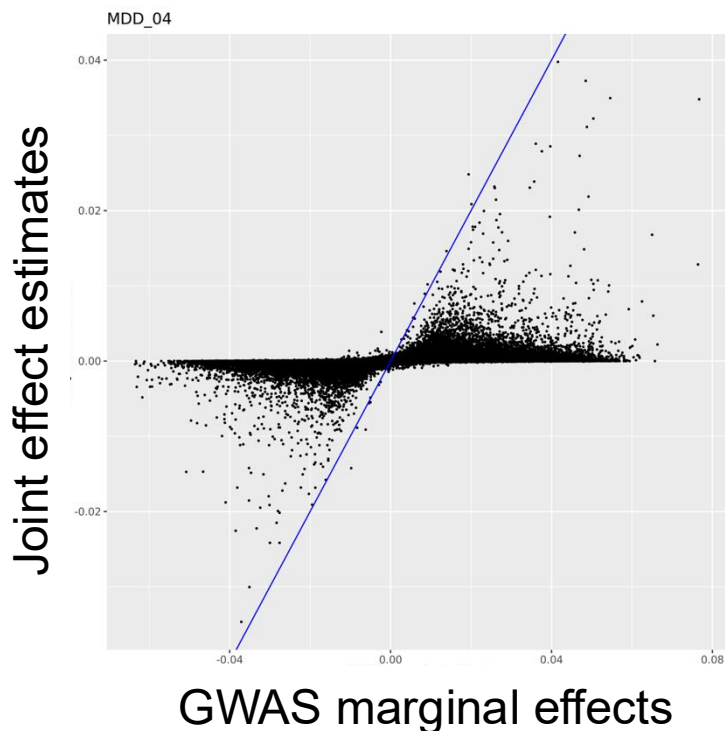
- No systematic differences in LD → **same ancestry**
- Minimum sampling variance in LD → LD ref sample size cannot be too small



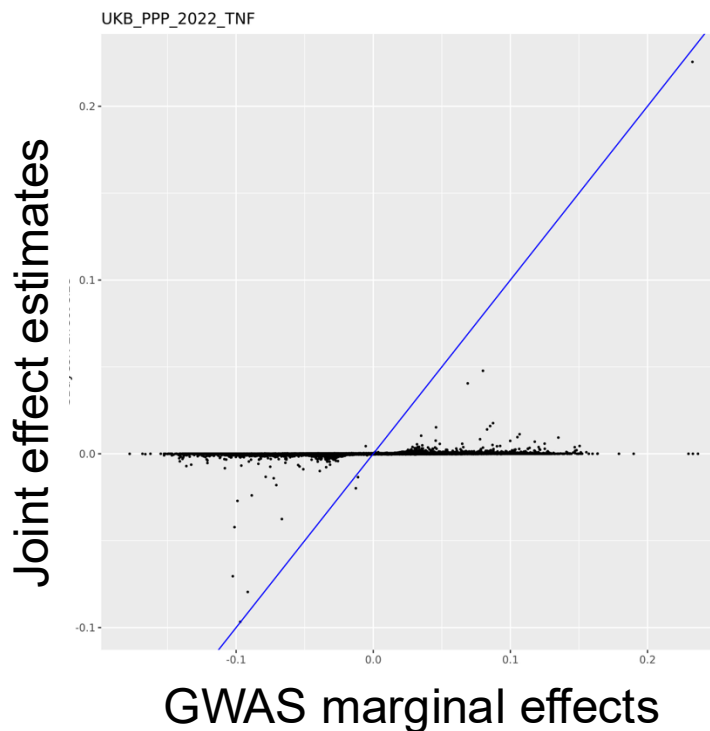
Failure to meet this assumption can result in a convergence issue!

GWAS marginal effect size vs. Estimated joint effect size

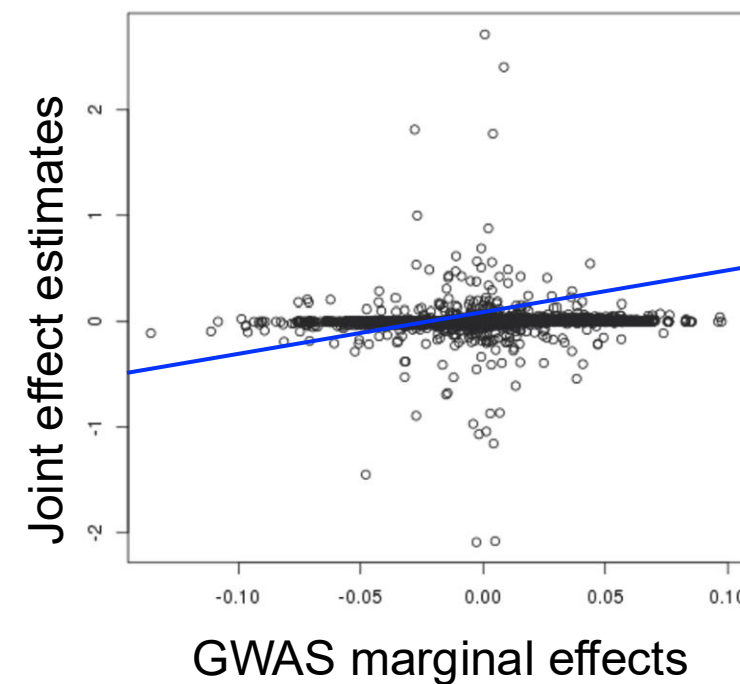
Most common 😊



Presence of large effects 😊



Bad convergence! 😞

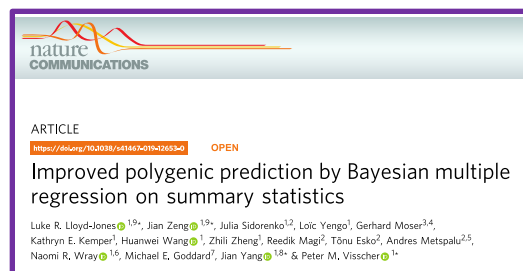


How do different methods handle this issue?

- Run multiple MCMC chains with different starting values
 - e.g., LDpred2
- Force an (ad hoc) minimum strong shrinkage to SNP effects
 - e.g., PRS-CS
- Regulate LD matrices
 - e.g., SBayesR uses chromosome-wide shrunk LD matrices
 - e.g., SBayesRC uses eigen-decomposed matrices from LD blocks

$$N \left(\frac{\mathbf{X}_j'(\mathbf{y} - \mathbf{1}_n\mu - \sum_{k \neq j} \mathbf{X}_k' \beta_k)}{\mathbf{X}_j' \mathbf{X}_j + \sigma_e^2 / \sigma_\beta^2}, \sigma_e^2 / (\mathbf{X}_j' \mathbf{X}_j + \sigma_e^2 / \sigma_\beta^2) \right)$$

the shrinkage parameter is set to be greater than a threshold value



Low-rank model (fits 7M SNPs or more)

In each quasi-independent LD block:

$$\mathbf{b} = \mathbf{R} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$

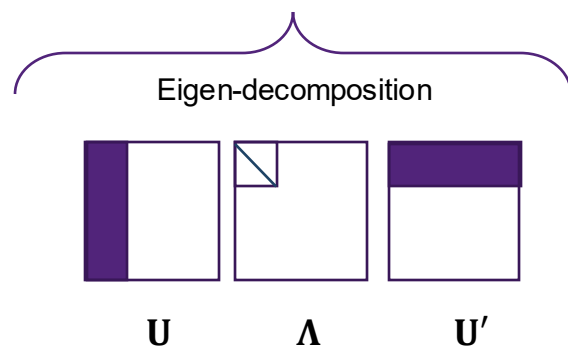
GWAS SNP marginal effects

LD correlation matrix

SNP joint effects

Residuals

$\text{Var}(\boldsymbol{\epsilon}) \propto$



$$\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}' \mathbf{b} = \boldsymbol{\Lambda}^{\frac{1}{2}} \mathbf{U}' \boldsymbol{\beta} + \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}' \boldsymbol{\epsilon}$$

$$\mathbf{w} = \mathbf{Q} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\text{Var}(\boldsymbol{\varepsilon}) \propto$

It only requires the top 20% PCs to explain 99.5% of the variance in LD!

Practical Exercise Part 2: SBayesR

<https://gctbhub.cloud.edu.au/data/teaching/>

Incorporating functional genomic annotations

Functional genomic annotations provide orthogonal information useful for polygenic prediction.

- Chromatin states
- Biological functions
- Molecular quantitative trait loci (xQTL)
-

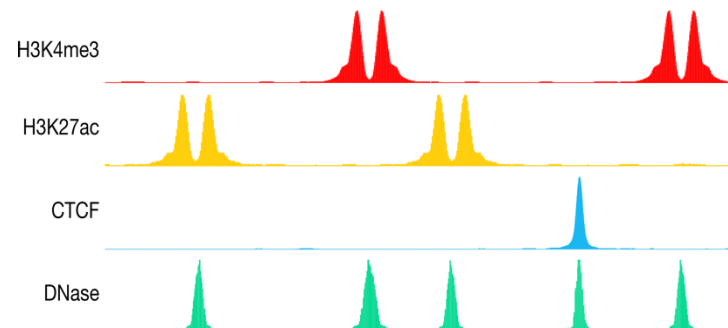
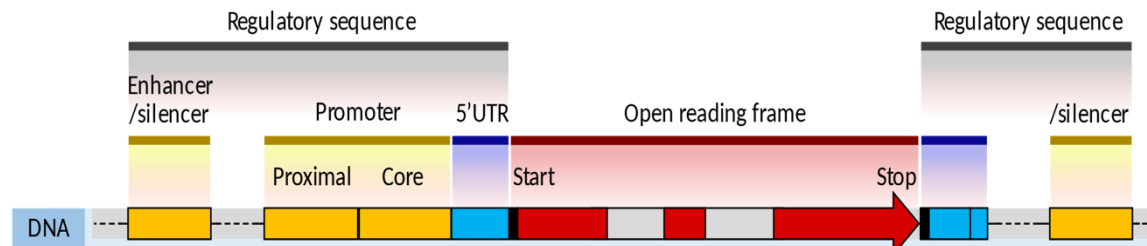
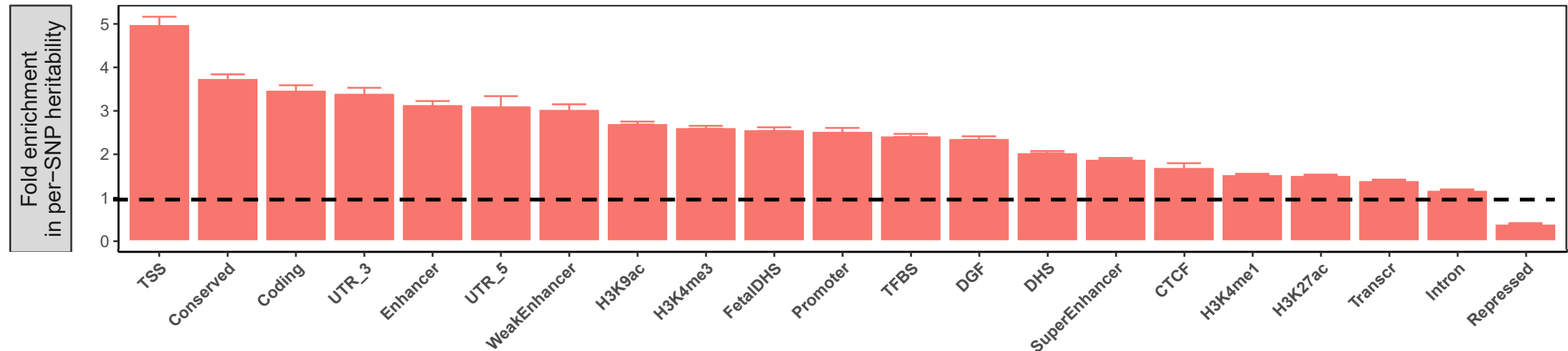


Image from ENCODE
CRICOS code 00025B

Functional genomic annotations provide orthogonal information useful for polygenic prediction.

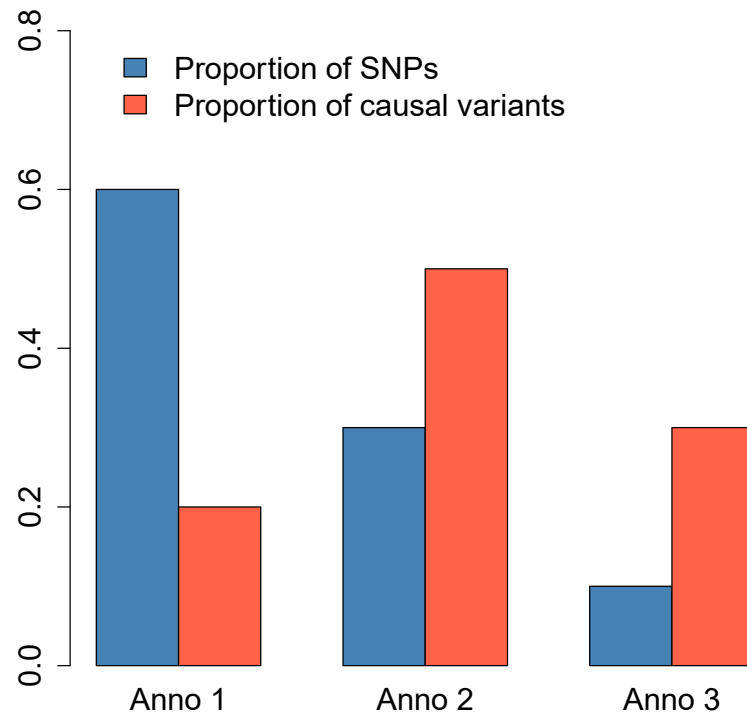
- Chromatin states
- Biological functions
- Molecular quantitative trait loci (xQTL)
-

Zeng et al 2021 Nature Communications

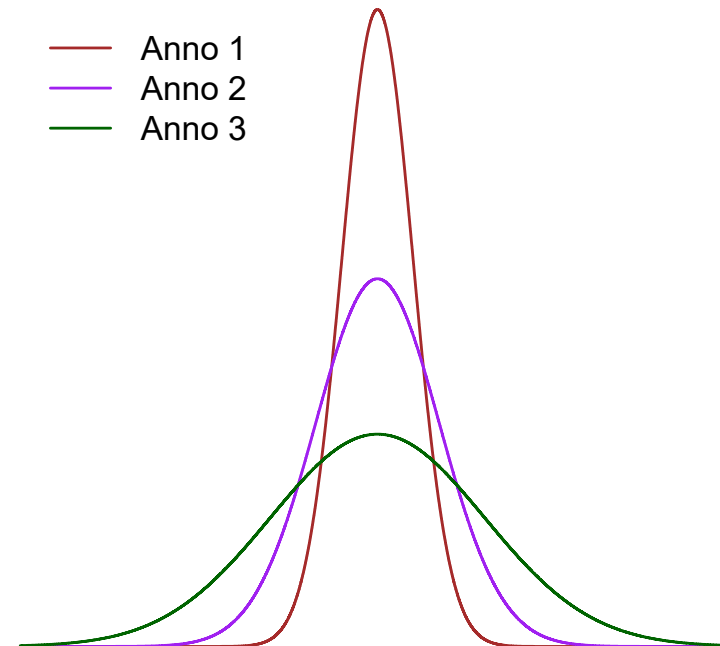


Functional annotations are informative on both the presence of causal variants and the distribution of causal effect sizes.

Differences in proportion of causal variants



Differences in distribution of causal effects



nature communications

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature communications](#) > [articles](#) > [article](#)

Article | [Open Access](#) | [Published: 18 October 2021](#)

Incorporating functional priors improves polygenic prediction accuracy in UK Biobank and 23andMe data sets

[Carla Márquez-Luna](#) , [Steven Gazal](#), [Po-Ru Loh](#), [Samuel S. Kim](#), [Nicholas Furlotte](#), [Adam Auton](#), [23andMe Research Team](#) & [Alkes L. Price](#) 

LDpred-funct

Exploiting biological priors and sequence variants enhances QTL discovery and genomic prediction of complex traits

[I. M. MacLeod](#) , [P. J. Bowman](#), [C. J. Vander Jagt](#), [M. Haile-Mariam](#), [K. E. Kemper](#), [A. J. Chamberlain](#), [C. Schrooten](#), [B. J. Hayes](#) & [M. E. Goddard](#)

[BMC Genomics](#) **17**, Article number: 144 (2016) | [Cite this article](#)

6209 Accesses | **146** Citations | **9** Altmetric | [Metrics](#)


BayesRC

PLOS COMPUTATIONAL BIOLOGY

 OPEN ACCESS  PEER-REVIEWED



RESEARCH ARTICLE

Leveraging functional annotations in genetic risk prediction for human complex diseases

[Yiming Hu](#) , [Qiongshi Lu](#) , [Ryan Powles](#), [Xinwei Yao](#), [Can Yang](#), [Fang Fang](#), [Xinran Xu](#), [Hongyu Zhao](#) 

AnnoPred

Winner's Curse Correction and Variable Thresholding Improve Performance of Polygenic Risk Modeling Based on Genome-Wide Association Study Summary-Level Data

[Jianxin Shi](#) , [Ju-Hyun Park](#), [Jubao Duan](#), [Sonja T. Berndt](#), [Winton Moy](#), [Kai Yu](#), [Lei Song](#), [William Wheeler](#), [Xing Hua](#), [Debra Silverman](#), [Montserrat Garcia-Closas](#), [Chao Agnes Hsiung](#), [Jonine D. Figueroa](#), [...], [Nilanjan Chatterjee](#)  [view all]

P+T-funct-LASSO

nature genetics

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature genetics](#) > [articles](#) > [article](#)

Article | [Published: 07 April 2022](#)

Leveraging fine-mapping and multipopulation training data to improve cross-population polygenic risk scores

[Omer Weissbrod](#) , [Masahiro Kanai](#), [Huwenbo Shi](#), [Steven Gazal](#), [Wouter J. Peyrot](#), [Amit V. Khera](#), [Yukinori Okada](#), [The Biobank Japan Project](#), [Alicia R. Martin](#), [Hilary K. Finucane](#) & [Alkes L. Price](#) 

[Nature Genetics](#) **54**, 450–458 (2022) | [Cite this article](#)

PolyPred

Need new method that can

- simultaneously fit all SNPs and annotation data in a unified model
- account for variations in both causal variant proportion and causal effect distribution

Leveraging functional annotations for cross-ancestry prediction

nature genetics



Article

<https://doi.org/10.1038/s41588-024-01704-y>

Leveraging functional genomic annotations and genome coverage to improve polygenic prediction of complex traits within and between ancestries

Received: 1 October 2022

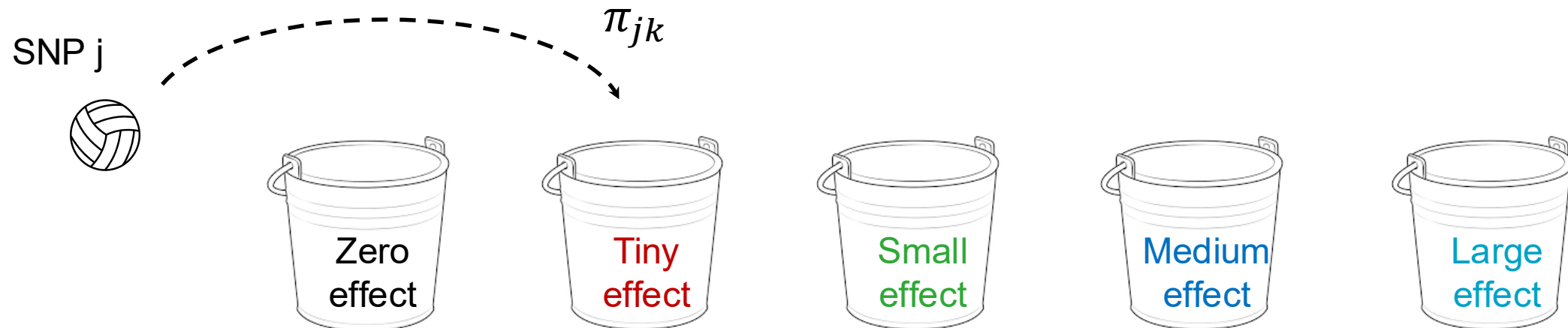
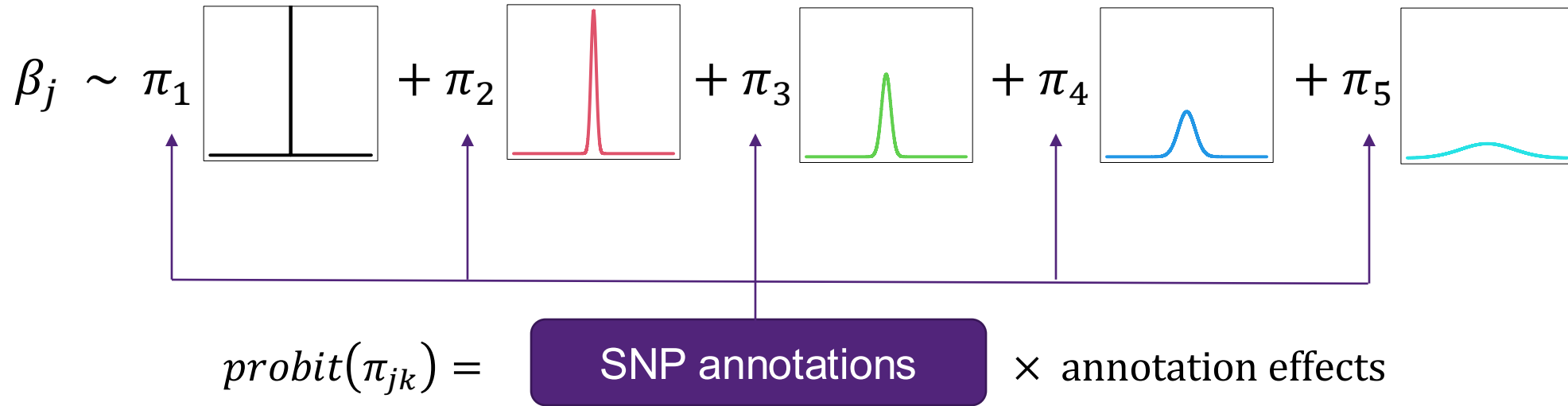
Accepted: 5 March 2024

Published online: 30 April 2024

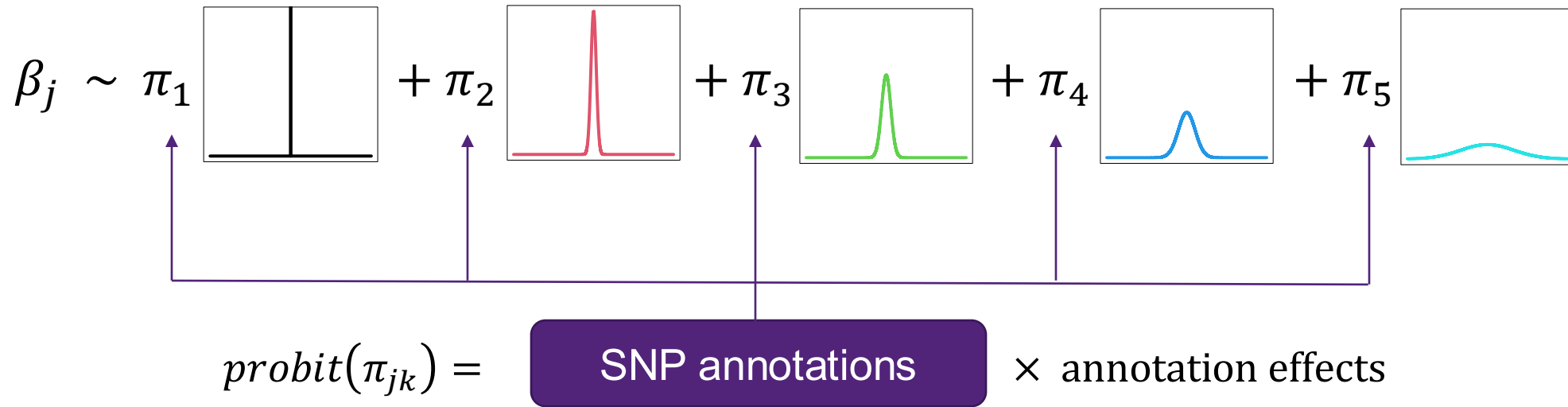
Check for updates

Zhili Zheng^{1,2,3}✉, Shouye Liu¹, Julia Sidorenko¹, Ying Wang¹, Tian Lin¹, Loic Yengo¹, Patrick Turley^{4,5}, Alireza Ani^{6,7}, Rujia Wang⁶, Ilja M. Nolte⁶, Harold Snieder⁶, LifeLines Cohort Study⁸, Jian Yang^{8,9}, Naomi R. Wray^{1,10}, Michael E. Goddard^{11,12}, Peter M. Visscher^{1,13} & Jian Zeng¹✉

Incorporate functional annotations through a hierarchical prior:



Incorporate functional annotations through a hierarchical prior:



Assumption

- Annotation effects are additive at the GLM scale.

Pros

- Estimation of conditional effects.
- Allow annotation overlap.
- Interpretation.

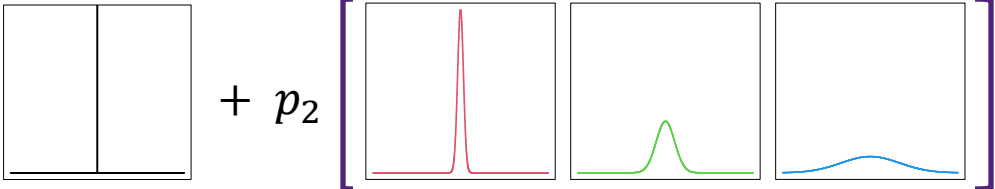
Cons

- # annotation effect parameters $\times 5$.
- $\pi_{j1} + \pi_{j2} + \pi_{j3} + \pi_{j4} + \pi_{j5} = 1$.

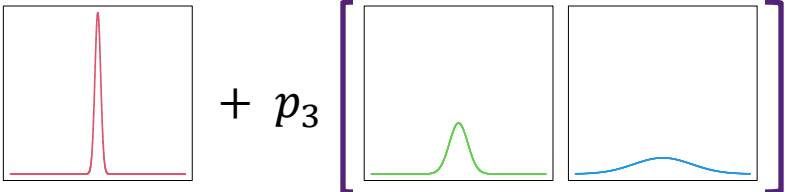
Suppose 4 components for simplicity

- A set of 2-component independent models:

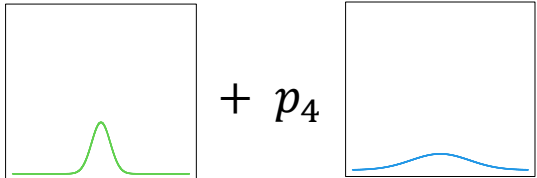
- For all SNPs

$$\beta_j \sim (1 - p_2) \left[\text{Null} \right] + p_2 \left[\text{Small} \mid \text{Medium} \mid \text{Large} \right]$$


- For SNPs with nonzero effects (conditional on non-null SNPs)

$$\beta_j \sim (1 - p_3) \left[\text{Small} \right] + p_3 \left[\text{Medium} \mid \text{Large} \right]$$


- For SNPs with at least medium effects (conditional on non-small-effect SNPs)

$$\beta_j \sim (1 - p_4) \left[\text{Medium} \right] + p_4 \left[\text{Large} \right]$$


p_2, p_3, p_4 are
independent!

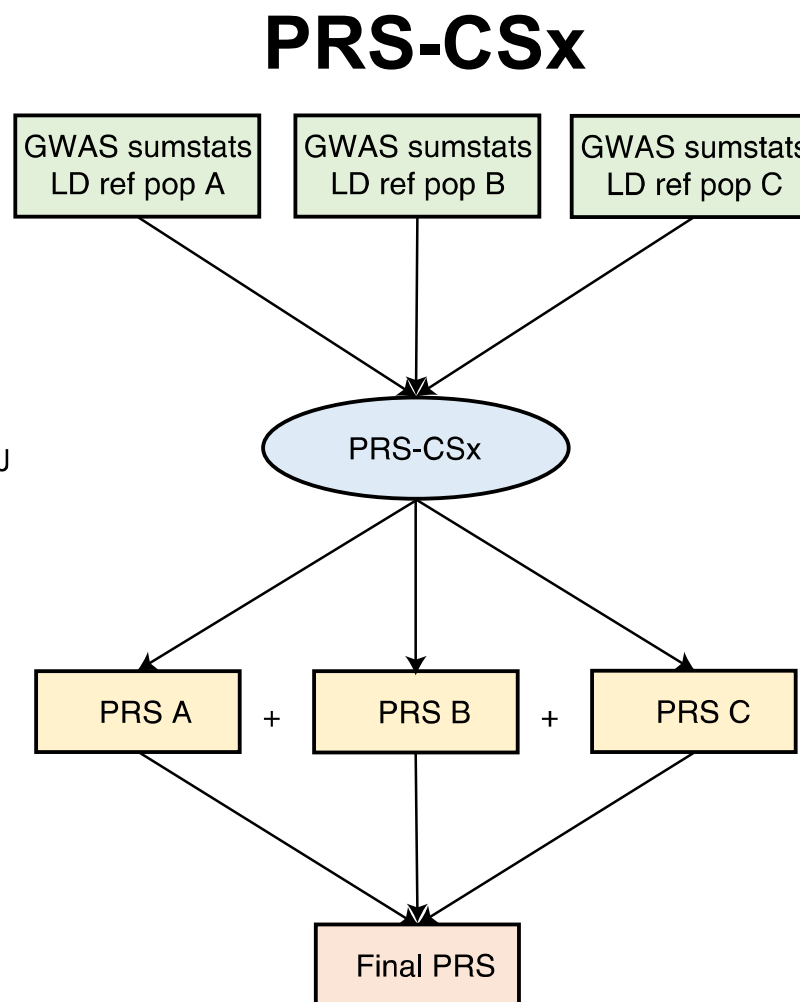
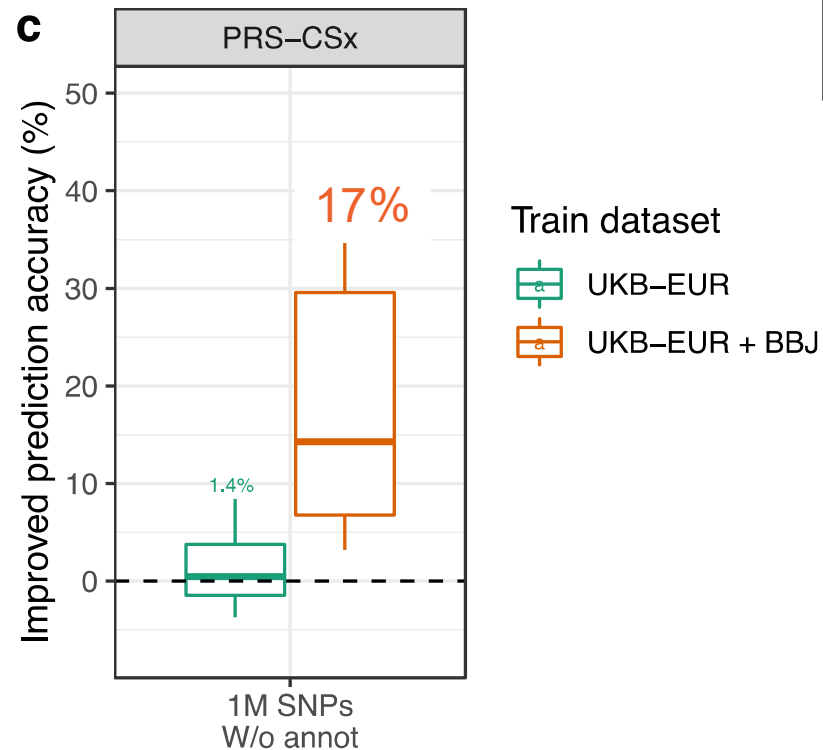
- Probit link function:

$$\Phi^{-1}(p) = \sum \text{SNP annotation} \times \text{annotation effect}$$

where Φ is the CDF of the standard normal distribution.

- It is straightforward to compute $p = \Phi(\cdot)$
and $\pi_1 = 1 - p_2$; $\pi_2 = (1 - p_3)p_2$; $\pi_3 = (1 - p_4)p_3p_2$; $\pi_4 = p_2p_3p_4$
- Assume a normal prior distribution for each annotation effect.
- Gibbs sampling for all parameters.

Use GWAS data from UKB EUR and BBJ EAS to predict UKB EAS



naturegenetics

Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

[nature](#) > [nature genetics](#) > [articles](#) > [article](#)

Article | Published: 05 May 2022

Improving polygenic prediction in ancestrally diverse populations

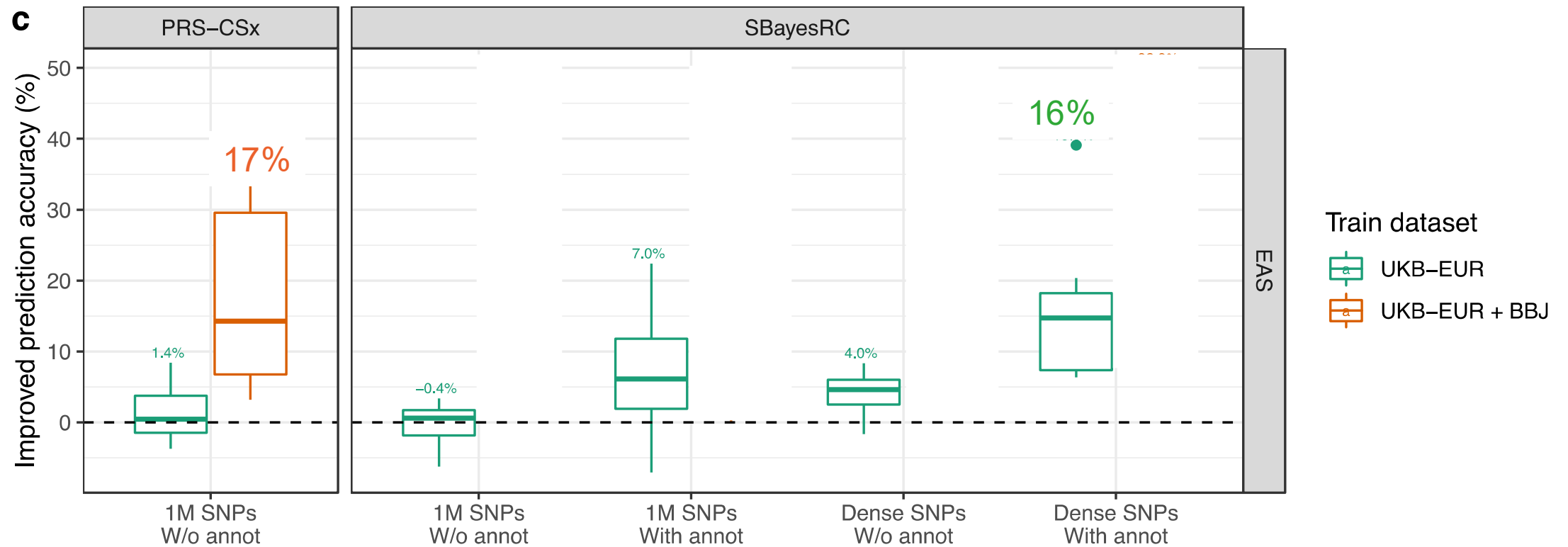
[Yunfeng Ruan](#), [Yen-Feng Lin](#), [Yen-Chen Anne Feng](#), [Chia-Yen Chen](#), [Max Lam](#), [Zhenglin Guo](#), [Stanley Global Asia Initiatives](#), [Lin He](#), [Akira Sawa](#), [Alicia R. Martin](#), [Shengying Qin](#) , [Hailiang Huang](#)  & [Tian Ge](#) 

Nature Genetics **54**, 573–580 (2022) | [Cite this article](#)

How important is functional annotation data compare to another GWAS dataset from the target ancestry?

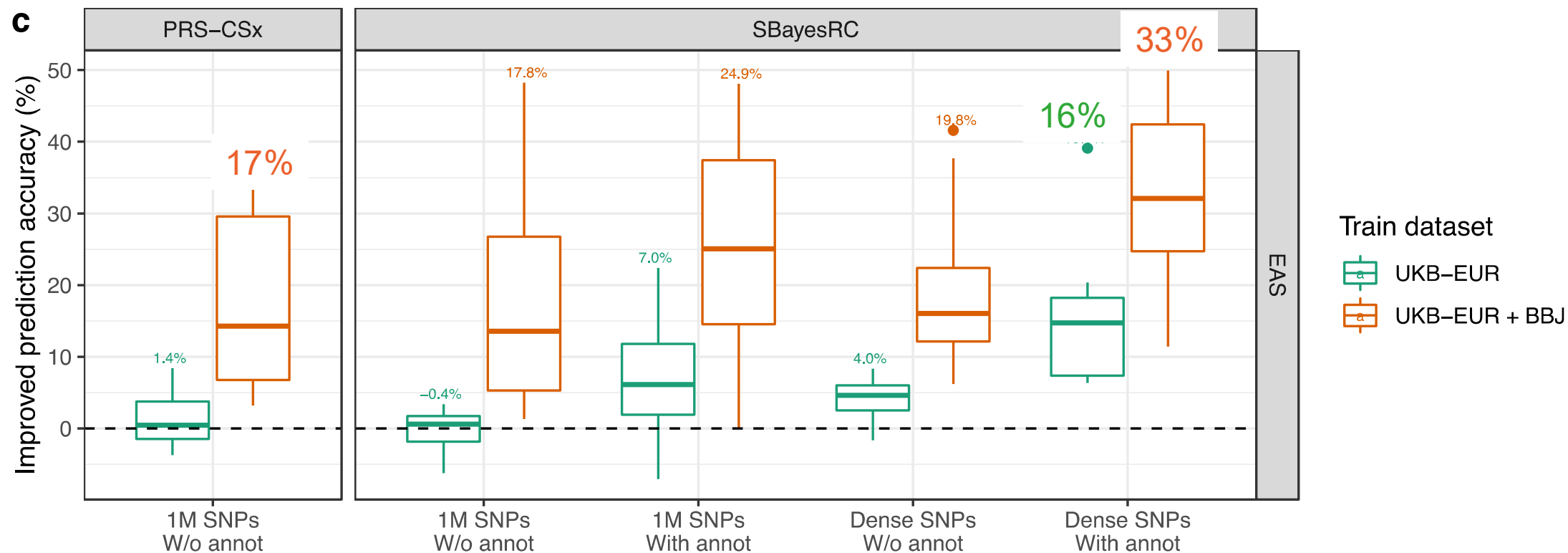
Trans-ancestry prediction

Use GWAS data from UKB EUR and BBJ EAS to predict UKB EAS



Trans-ancestry prediction

Use GWAS data from UKB EUR and BBJ EAS to predict UKB EAS



Practical Exercise Part 3: SBayesRC

<https://gctbhub.cloud.edu.au/data/teaching/>

Summary and challenges

- PGS are imperfect genetic predictors with inherently limited accuracy.
- Being improved with more data and better methods.
- Bayesian approach allows us to incorporate prior knowledge in estimation of SNP effects.
- State-of-the-art Bayesian methods utilize GWAS summary statistics, which unleash the power of large GWAS sample size but may encounter convergence issue.
- Different methods mainly differ in prior assumption and implementation.
- Functional annotations with genome coverage provide orthogonal information to GWAS data and therefore boost prediction accuracy.

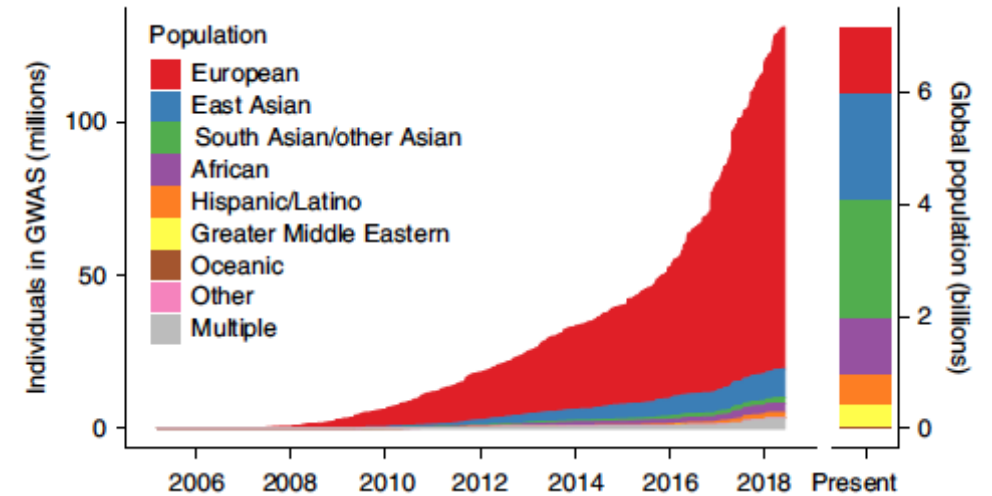
PERSPECTIVE

<https://doi.org/10.1038/s41588-019-0379-x>

nature
genetics

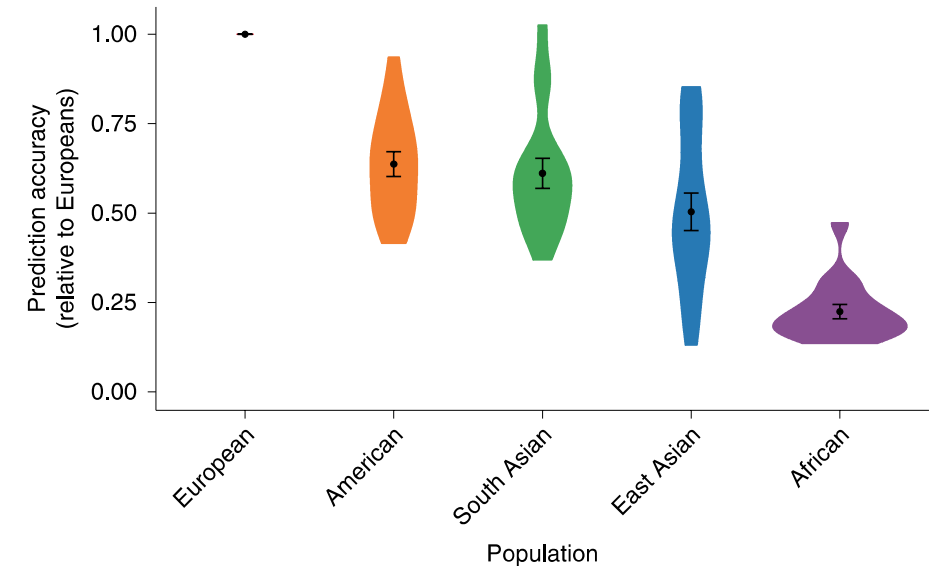
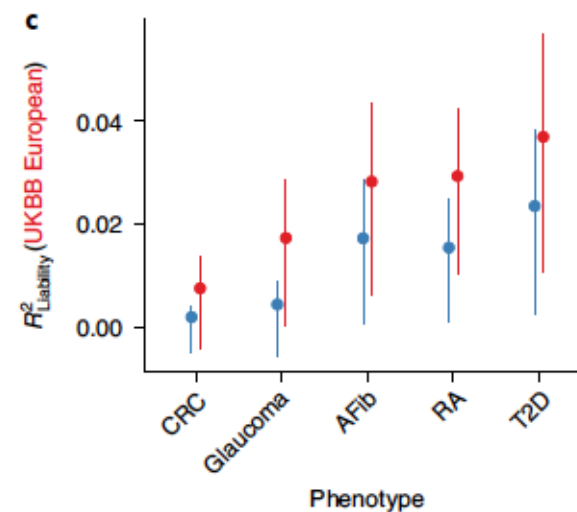
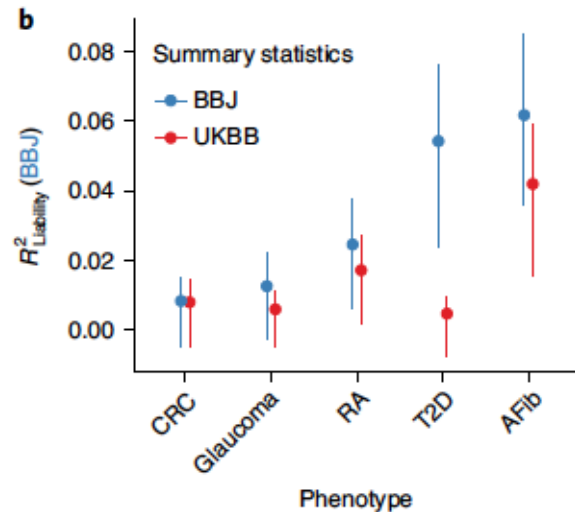
Clinical use of current polygenic risk scores may exacerbate health disparities

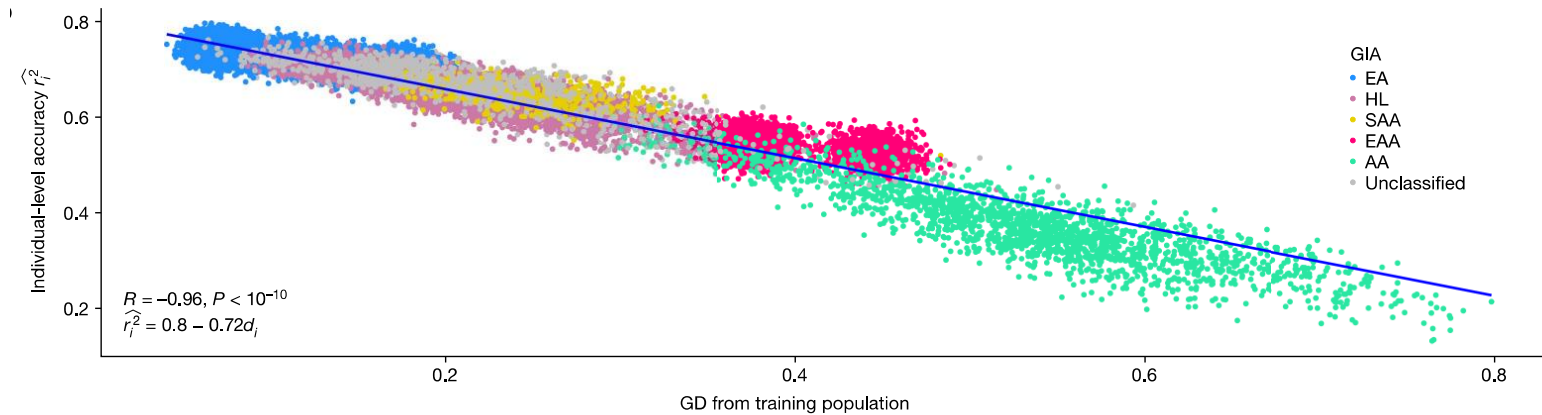
Alicia R. Martin ^{1,2,3*}, Masahiro Kanai ^{1,2,3,4,5}, Yoichiro Kamatani ^{5,6}, Yukinori Okada ^{5,7,8}, Benjamin M. Neale ^{1,2,3} and Mark J. Daly ^{1,2,3,9}



Predicted into Japanese

Predicted into European





nature

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [articles](#) > article

Article | [Open Access](#) | Published: 17 May 2023

Polygenic scoring accuracy varies across the genetic ancestry continuum

Yi Ding , Kangcheng Hou, Zigi Xu, Aditya Pimplaskar, Ella Petter, Kristin Boulier, Florian Privé, Bjarni J. Vilhjálmsson, Loes M. Olde Loohuis & Bogdan Pasaniuc

Nature (2023) | [Cite this article](#)

11k Accesses | 200 Altmetric | [Metrics](#)

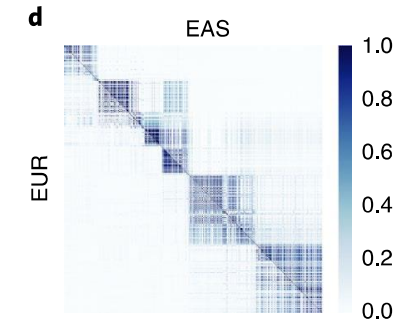
Issues

- Same causal variants
 - Different allele frequencies
 - LD differences
 - Different effect sizes
- Different causal variants
 - GxE
 - Different phenotype

In general:

We expect common causal variants to be shared across ancestries

But correlation structure differs



nature genetics

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [nature genetics](#) > [articles](#) > article

Article | Published: 20 March 2023

Causal effects on complex traits are similar for common variants across segments of different continental ancestries within admixed individuals

nature genetics

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [nature genetics](#) > [articles](#) > article

Article | [Open access](#) | Published: 03 February 2025

Fine-scale population structure and widespread conservation of genetic effect sizes between human groups across traits

What is the maximum prediction accuracy we can get?

Variance explained by
the predictor

$$R^2 = \frac{h_m^2}{1 + C}$$

h_m^2 : **True** variance explained by the predictor
depends on the SNP set - subscript m.

C: captures the error in estimation

As $C \rightarrow 0$, $R^2 \rightarrow h_m^2$

We want C to be as small as possible:

- C decreases as Discovery sample **N** **increases**
- C decreases as the number of SNPs in the SNP set **m** **decreases**

$$C \approx \frac{m}{Nh_m^2}$$



As m gets smaller, h_m^2 also gets smaller

How to optimise m and h_m^2 to get max R^2 ?



How about whole genome sequencing?

Maximum depends on
maximising h_m^2

We use GWAS data so the
maximum h_m^2 is the SNP-based
heritability

Theoretical maximum depends
on the heritability of the trait

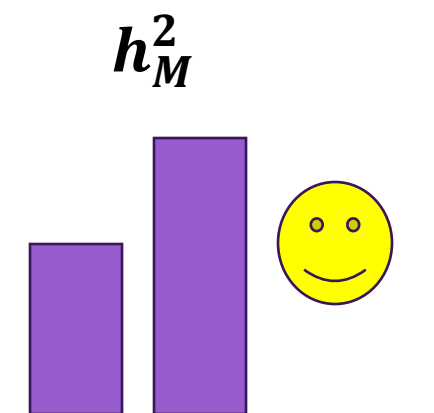
$$R^2 \approx \frac{h_m^2}{1 + \frac{m}{Nh_m^2}}$$

With whole genome sequencing the variance captured by
all measured SNPs will increase

But the number of SNPs that we have estimate effect sizes
for increases much more

Need MASSIVE discovery sample sizes for WGS associations

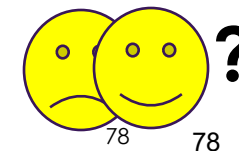
Also... rare variants are less likely to be shared across populations



M



R^2



Reference

BayesA, BayesB:

Copyright © 2001 by the Genetics Society of America

Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps

T. H. E. Meuwissen,* B. J. Hayes† and M. E. Goddard†‡

*Research Institute of Animal Science and Health, 8200 AB Lelystad, The Netherlands, †Victorian Institute of Animal Science, Attwood 3049, Victoria, Australia and ‡Institute of Land and Food Resources, University of Melbourne, Parkville 3052, Victoria, Australia

Manuscript received August 17, 2000
Accepted for publication January 17, 2001

BayesC:

Habier et al. *BMC Bioinformatics* 2011, **12**:186
<http://www.biomedcentral.com/1471-2105/12/186>



RESEARCH ARTICLE

Open Access

Extension of the bayesian alphabet for genomic selection

David Habier^{1*}, Rohan L. Fernando¹, Kadir Kizilkaya^{1,2} and Dorian J. Garrick^{2,3}

BayesR:



J. Dairy Sci. 95:4114–4129

<http://dx.doi.org/10.3168/jds.2011-5019>

© American Dairy Science Association®, 2012. Open access under [CC BY-NC-ND license](#)

Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels

M. Erbe,*¹ B. J. Hayes,†‡^{1,2} L. K. Matukumalli,§ S. Goswami,|| P. J. Bowman,†‡ C. M. Reich,†‡ B. A. Mason,†‡ and M. E. Goddard†‡¶



RESEARCH ARTICLE

Simultaneous Discovery, Estimation and Prediction Analysis of Complex Traits Using a Bayesian Mixture Model

Gerhard Moser^{1*}, Sang Hong Lee¹, Ben J. Hayes^{2,3}, Michael E. Goddard^{2,4}, Naomi R. Wray¹, Peter M. Visscher^{1,5}

SBayesR:



ARTICLE

<https://doi.org/10.1038/s41467-020-18653-0> OPEN

Improved polygenic prediction by Bayesian multiple regression on summary statistics

Luke R. Lloyd-Jones^{1,9*}, Jian Zeng^{1,9*}, Julia Sidorenko^{1,2}, Loic Yengo¹, Gerhard Moser^{2,4}, Kathryn E. Kemper¹, Huanwei Wang¹, Zhili Zheng¹, Reedik Magi², Tõnu Esko², Andres Metspalu^{2,5}, Naomi R. Wray^{1,6}, Michael E. Goddard⁷, Jian Yang^{1,8*} & Peter M. Visscher^{1*}

SBayesRC:

nature genetics

Article

<https://doi.org/10.1038/s41588-024-01704-y>

Leveraging functional genomic annotations and genome coverage to improve polygenic prediction of complex traits within and between ancestries

Received: 1 October 2022

Accepted: 5 March 2024

Published online: 30 April 2024

Check for updates

Zhili Zheng^{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000}

LDpred2:

Bioinformatics, 2021, 1–8
doi: 10.1093/bioinformatics/btaa1029
Advance Access Publication Date: 16 December 2020
Original Paper



Genetics and population analysis

LDpred2: better, faster, stronger

Florian Privé^{1,*}, Jülyan Arbel² and Bjarni J. Vilhjálmsson^{1,3,*}

PRS-CS:

nature communications

Explore content ▾ About the journal ▾ Publish with us ▾

nature > [nature communications](#) > [articles](#) > article

Article | [Open access](#) | Published: 16 April 2019

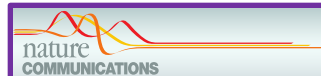
Polygenic prediction via Bayesian regression and continuous shrinkage priors

Tian Ge¹, Chia-Yen Chen¹, Yang Ni¹, Yen-Chen Anne Feng¹ & Jordan W. Smoller¹

Nature Communications 10, Article number: 1776 (2019) | [Cite this article](#)

70k Accesses | 1580 Citations | 53 Altmetric | [Metrics](#)

MegaPRS:



ARTICLE

<https://doi.org/10.1038/s41467-021-24485-y> OPEN

Improved genetic prediction of complex traits from individual-level data or summary statistics

Qianqian Zhang¹, Florian Privé², Bjarni Vilhjálmsson^{1,2} & Doug Speed^{1,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500}