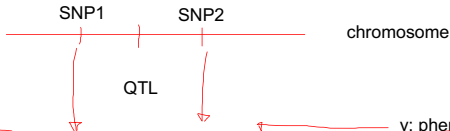


Genomic BLUP

Peter von Rohr

2021-03-08

So Far



a1: marker effect for SNP1

- ▶ Estimate effect of few SNP loci linked to QTL
- ▶ Use parameter estimates to predict genomic breeding values
- ▶ **New:** Many SNP, find the important ones

150'000 - 20 Mio

denoted by the variable p

combine estimation of marker effects together with availability of many SNP in the dataset

y: phenotypes

N

$\}$

$y =$

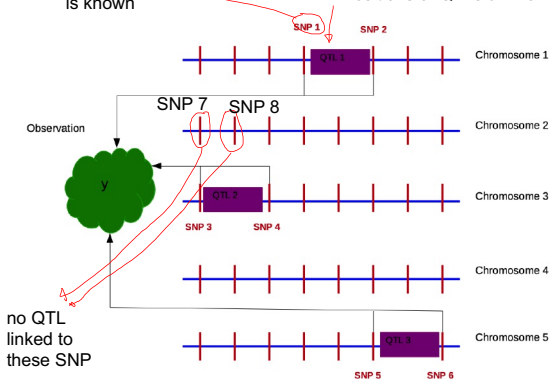
$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

$$P \gg N$$

least squares cannot be used to estimate marker effects

Situation SNP positions
is known

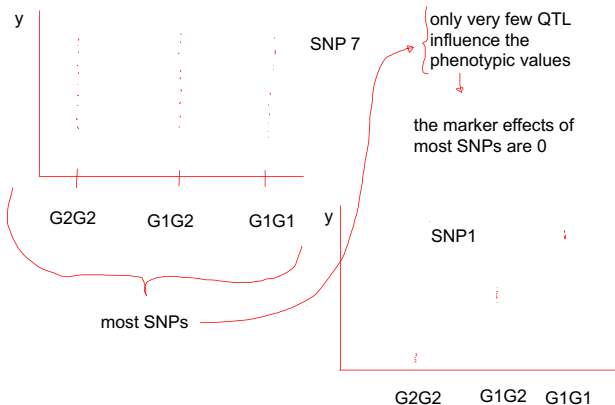
Positions of QTLs unknown



no QTL
linked to
these SNP

Goal: Find SNP 1 – SNP 6 out of the many SNPs

For the two SNPs 7 and 8 that are not linked to a QTL



Summary: 2 Problems

1. if we consider all SNP in our data set, then $p \gg N$
==> least squares cannot be used
2. from genetic model: only few QTL for a given trait,
==> most SNP have marker effects (a) = 0

because the position of the QTL is unknown, we do not know which SNP have marker effects = 0

Approaches in Fixed Linear Model Framework

Possible solution for problem 2: Model selection to determine which SNP have marker effect that are not 0


Two Approaches

1. Forward selection: Start with empty model, include predictors that improve model
2. Backward elimination: Start with full model, remove predictors as long as model does not get worse

no SNP



all SNP, not possible
due to $p \gg N$



Forward Selection

Start with empty model

no SNP

intercept

$$y = b_0 + e$$

$SSR = r_1^2 + r_2^2 + \dots + r_n^2$

Include additional predictor

$M_1: y = b_0 + b_1 + e$

$M_2: y = b_0 + b_1 + b_2 + e$

$M_p:$

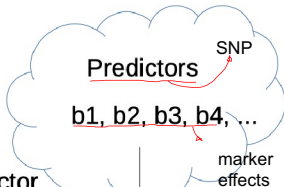
choose b_1 such that new SSR is reduced the most

Model better ?

No, ignore b_1

Criteria: Mallow Cp, Akaike

Yes, keep b_1



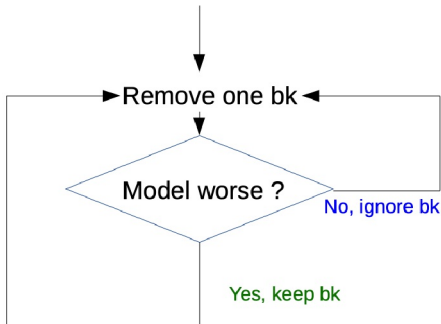
Backward Elimination

except for SNP data, this is the preferred way

Start with full model

not possible with
SNP, because $p \gg N$

$$Y = b_0 + b_1 + b_2 + b_3 + \dots + e$$



Model Selection With Genomic Data

cannot use backward elimination with genomic data, because parameter estimation in the full model cannot be

- ▶ Only backward elimination really works in practical problems
- ▶ Large number of predictors ($1.5 * 10^5$)
- ▶ How to determine sequence of predictors to eliminate
- ▶ Fitting the full model is problematic

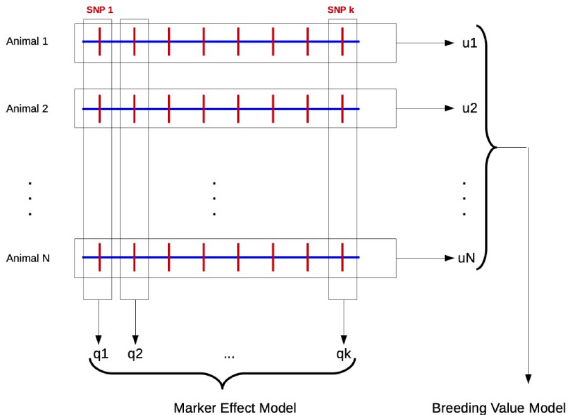
Mixed Linear Effect Model

- ▶ One solution: replace fixed linear effect model by mixed linear effect model (mle)
- ▶ MLE: additional random effect besides error term
- ▶ Random effects are specified by expected value and variance
- ▶ In livestock breeding MLE have a good reputation from BLUP animal model

MLE In Genomics

- ▶ Two different parametrizations
 1. Marker Effect Model (MEM)
 2. Breeding Value Model (BVM)

Overview



Marker Effect Model

In MEM random effects of markers are directly included in the model. For an idealized data set we can write

$$y = 1_n\mu + Wq + e$$

where

- y vector of length n with observations
- μ general mean denoting fixed effects
- 1_n vector of length n of all ones
- q vector of length m of random SNP effects
- W design matrix relating SNP-genotypes to observations
- e vector of length n of random error terms

Breeding Value Model

$$y = Xb + Zg + e$$

where

- y vector of length n with observations
- b vector of length r with fixed effects
- X incidence matrix linking elements in b to observations
- g vector of length t with random genomic breeding values
- Z incidence matrix linking elements in g to observations
- e vector of length n of random error terms