

SISG  
2014

AGTGAAGCTACCTACTTAGAAAGTGACTGCTACTGGTGAAAAT

## SISG Module 23: Advanced Quantitative Genetics

19th Summer Institute in Statistical Genetics

**W** UNIVERSITY *of* WASHINGTON

(This page left intentionally blank.)

## Overview/reminder of basic concepts in statistics and genetics

Matt Robinson, University of Queensland

1

### Bayes' Theorem

Identify people who are liable to suffer from a genetic disease later in life.

1 in 1000 people are a carrier of the disease

No test is perfect - probability that a carrier tests negative is 1%  
- probability that a non-carrier tests positive is 5%

A = the event “the patient is a carrier”

B = the event “the test result is positive”

Hence:  $P(A) = 0.001$ ;  $P(A') = 0.999$ ;  $P(B|A) = 0.99$ ;  $P(B|A') = 0.05$

A patient has a positive result. Q: What is the probability that the patient is a carrier?

#### Answer

$$\begin{aligned} P(A|B) &= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A')P(A')} \\ &= \frac{0.99 * 0.001}{(0.99 * 0.001) + (0.05 * 0.999)} = 0.0194 \end{aligned}$$

2

## Random variables, expected values and (co)variance

A discrete random variable can assume only a countable number of values

Probability mass function:

$$p(x) = P(X = x)$$

Expected value:

$$\mu = E(X) = \sum xp(x)$$

As a function of random variable:

$$E[h(X)] = \sum h(x)p(x)$$

Variance:

$$Var(X) = E[(X - \mu)^2]$$

3

## Random variables, expected values and (co)variance

A discrete random variable can assume only a countable number of values

allele x:

$$p(x) = \begin{cases} p & x = 1 \\ 1-p & x = 0 \end{cases}$$

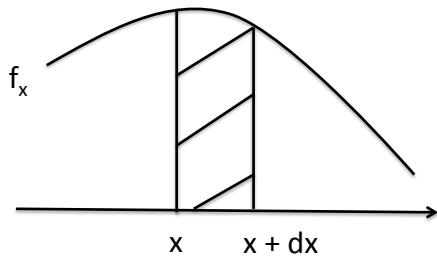
$$E(X) = 0(1-p) + 1(p) = p$$

$$Var(X) = p - p^2 = p(1-p)$$

4

## Random variables, expected values and (co)variance

A continuous random variable can be any value within a range



probability of being in shaded area

$$= f_X(x)dx$$

the interval should contribute

$$= xf_X(x)dx$$

the expected value and variance

$$E(X) = \mu_X = \int_{-\infty}^{\infty} xf_X(x)dx \quad Var(X) = E((X - \mu_X)^2)$$

5

## Random variables, expected values and (co)variance

Covariance

Let X and Y be a pair of continuous random variables, with respective means  $\mu_X$  and  $\mu_Y$ . The expected value of  $(X - \mu_X)(Y - \mu_Y)$  is called the covariance between X and Y.

$$Cov(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

If the random variables X and Y are independent, then the covariance between them is 0. However, the converse is not true.

6

## Summary (co)variance rules

$$Var(x) = E[x - E(x)]^2$$

$$Var(cx) = c^2 Var(x)$$

$$Var(x + y) = Var(x) + Var(y) + 2Cov(x, y)$$

$$Var(x + c) = Var(x)$$

$$Cov(x, y) = E[(x - E(x))(y - E(y))]$$

$$Cov(cx, y) = cCov(x, y)$$

$$Cov(x, y + z) = Cov(x, y) + Cov(x, z)$$

7

## Hardy-Weinberg equilibrium

Mathematical relation between **allele** frequencies and the **genotype** frequencies is:

		Allele Frequency	A	a
Allele	Frequency	p	$p^2$	$2pq$
		q	$qp$	$q^2$

8

## HWE and SNPs

If SNP genotypes are coded  $x = 0, 1$  and  $2$  (alleles) and the allele frequency is  $p$ , then:

$$Var(x) = 2p(1-p)$$

# Resemblance between relatives

Mike Goddard

1

What do we mean by resemble?

Similar values of quantitative traits

Measure by correlation

$$= \text{Covariance}(y_i, y_j) / \text{variance}(y)$$

2

Why do relatives resemble each other?

3

Why do relatives resemble each other?

Similar

Genes

Family environment

Country

School

4

# Model phenotype

Phenotype = genetic effect

- + country
- + year of birth
- + family environment

Fixed effects

Country, year of birth

Random effects

Genetic effect, family environment

We need a model of the covariances between terms

5

# Model phenotype

Phenotype = genetic effect

- + country
- + year of birth
- + family environment
- + individual environment

$$V(\text{phenotype}) = V(\text{genetic effects}) + V(\text{family environment}) + V(\text{individual environment})$$

$$\text{Cov}(\text{phenotype}_i, \text{phenotype}_j) = \text{Cov}(\text{genetic effects}) + \text{Cov}(\text{family environments})$$

6

## Model phenotype

Random effects

Genetic effect, family environment

We need a model of the covariances between terms

$$C(\text{family environments}) = \begin{cases} 0 & \text{if different families} \\ 1 * V_{CE} & \text{if same family} \end{cases}$$

7

## Covariance between genetic effects of relatives

Model with 1 gene, 2 alleles and additive gene action

We need genetic variances and covariances

Genotype	BB	Bb	bb
Effect	a	0	-a
Frequency	$p^2$	$2pq$	$q^2$

$$(p+q=1)$$

$$\text{Mean} = a * p^2 + 0 * 2pq - a * q^2 = (p-q)*a$$

$$\begin{aligned}\text{Variance (genetic effect)} &= \text{genetic variance} = V_G \\ &= E(\text{effect}^2) - E(\text{effect})^2\end{aligned}$$

$$V_G = a^2 * p^2 + 0 * 2pq + a^2 * q^2 - [(p-q)*a]^2 = 2pqa^2$$

8

## Model with 1 gene, 2 alleles and additive gene action

Covariance between parent and offspring

Parent		Offspring			
Genotype	frequency	BB	Bb	bb	mean
BB	a	$p^2$	p	q	$pa$
Bb	0	$2pq$	$0.5p$	0.5	$0.5q$
bb	-a	$q^2$		p	$-qa$

$$\text{Cov}(\text{parent genetic value, offspring genetic value}) \\ = p^2 * a * pa + q^2 * (-a) * (-qa) - [(p-q)a] * [(p-q)a] = pqa^2 = 0.5 V_G$$

9

## Model with 1 gene, 2 alleles and additive gene action

Covariance between parent and offspring (another way)

Model genetic value as sum of gametic effects from mother and father

$$g = x_m + x_f$$

$$V(g) = V(x_m) + V(x_f) = 2V(x)$$

$$C(g_p, g_o) = C(x_{mp} + x_{fp}, x_{mo} + x_{fo}) \\ = C(x_{mp}, x_{mo}) + C(x_{mp}, x_{fo}) + C(x_{fp}, x_{mo}) + C(x_{fp}, x_{fo}) \\ = 0 \quad \quad \quad + ? \quad \quad \quad + 0 \quad \quad + ?$$

$$C(x_{mp}, x_{fo}) = V(x) \text{ if } x_{mp} \text{ is ibd to } x_{fo} \\ = 0 \text{ otherwise}$$

$$C(x_{mp}, x_{fo}) = C(x_{fp}, x_{fo}) = 0.5 V(x)$$

$$C(g_p, g_o) = 0 + 0.5V(x) + 0 + 0.5V(x) = V(x) = 0.5 V_G$$

10

## Probability that relatives share alleles IBD

Covariance between relatives depends on probability that their alleles are IBD

This probability can be calculated from pedigrees

Assume that base individuals at the top of the pedigree (ie those without a pedigree) have unrelated alleles ie the individuals are unrelated

Recurrence formulae for P(IBD)

if i and j are base individuals,  $P(x_{.i} \equiv x_{.j}) = 0$

Otherwise,  $P(x_{.i} \equiv x_{fj}) = 0.5 [P(x_{.i} \equiv x_{fk}) + P(x_{.i} \equiv x_{mk})]$  where k is the father of j

11

## Probability that relatives share alleles IBD

$$\begin{matrix} & k (mk, fk) \\ & \diagdown \\ i (mi, fi) & j(mj, fj) \end{matrix}$$

12

## Relationships between individuals

$P(\text{gametes are IBD})$  can be stored in a gametic relationship matrix  
 $G(w_i, z_j) = P(w_i \equiv z_j)$

But usually we analyse measurements on diploid individuals

$$C(g_i, g_j) = A(i, j) V_G = [G(m_i, m_j) + G(m_i, f_j) + G(f_i, m_j) + G(f_i, f_j)] V(x) \\ = [G(m_i, m_j) + G(m_i, f_j) + G(f_i, m_j) + G(f_i, f_j)] V_G / 2$$

$$A(i, j) = [G(m_i, m_j) + G(m_i, f_j) + G(f_i, m_j) + G(f_i, f_j)] / 2$$

where  $A$  is the numerator relationship matrix

13

## Relationships between individuals

Example: Relationship of individual with herself

Gametic relationship matrix

	mi	fi
mi	1	0
fi	0	1

$$\text{Numerator relationship } A(i, i) = [1+0+0+1]/2 = 1$$

14

## Relationships between individuals

Example: Relationship of sisters

Gametic relationship matrix

	mi	fi
mj	0.5	0
fj	0	0.5

Numerator relationship  $A(i,j) = [0.5+0+0+0.5]/2 = 0.5$

15

## Relationships between individuals

$i = (im, if)$  and  $j = (jm, jf)$

Co-ancestry of i and j

= Inbreeding co-efficient of an offspring of i and j

= Prob( offspring gets two alleles that are IBD)

=  $(P(im \equiv jm) + P(im \equiv jf) + P(if \equiv jm) + P(if \equiv jf))/4$

=  $A(i,j) / 2$

Additive relationship (NRM) =  $2 * \text{co-ancestry}$

=  $2 * \text{kinship}$

16

## Estimating genetic variance

Data on phenotypes ( $y$ ) of related subjects

$$y = \text{fixed effects} = g + e$$

$$V(g) = A V_G$$

$$V(e) = I V_E$$

Use ML or REML to estimate variances

17

## Estimating genetic variance

Use ML or REML to estimate variances

ML finds the value of  $V_G$  that maximises the probability of observing the data

ML estimates all parameters together

= estimates variances assuming that fixed effects have been estimated without error

REML allows for loss of df in estimating fixed effects

ML  $\sigma^2 = \sum(y\text{-mean})^2/N$

REML  $\sigma^2 = \sum(y\text{-mean})^2/(N-1)$

Little difference unless many fixed effects

Use REML computer programs such as ASREML

18

## Estimating genetic variance

Example: Data on phenotypes ( $y$ ) of full sibs

$y = \text{fixed effects} = g + e$

$\text{Cov}(g_i, g_j) = A(i,j)$   $V_G = 0.5 V_G$  if  $i$  and  $j$  are sibs

Therefore estimate  $V_G$  by  $2\text{cov}(\text{full-sibs})$

$h^2$  by  $2$  correlation between full-sibs

What is the covariance between twins?

19

## Model with dominance

20

## Covariance between genetic effects of relatives

Model with 1 gene, 2 alleles and additive and dominant gene action  
We need genetic variances and covariances

Genotype	BB	Bb	bb	
Effect	a	d	-a	
Frequency	$p^2$	$2pq$	$q^2$	$(p+q=1)$

$$\text{Mean} = a * p^2 + d * 2pq - a * q^2 = (p-q)*a + 2pqd$$

$$\text{Variance (genetic effect)} = \text{genetic variance} = V_G$$

$$= E(\text{effect}^2) - E(\text{effect})^2$$

$$V_G = a^2 * p^2 + d^2 * 2pq + a^2 * q^2 - [(p-q)*a + 2pqd]^2 = 2pq\alpha^2 + (2pqd)^2$$

$$\text{where } \alpha = a + (q-p)d$$

21

## Covariance between genetic effects of relatives

Model with 1 gene, 2 alleles and additive and dominant gene action  
but the covariance between relatives doesn't depend directly on VG. We need to decompose VG into an additive and dominance variance.

Parameterise the genetic value as

$$g = \text{mean} + \text{additive effect} + \text{dominance deviation}$$

$$g = \text{mean} + \text{paternal allele effect} + \text{maternal allele effect} + \text{interaction of alleles}$$

Genotype	BB	Bb	bb	
Effect	a	d	-a	
Frequency	$p^2$	$2pq$	$q^2$	$(p+q=1)$
mean	$(p-q)a + 2pqd$	$(p-q)a + 2pqd$	$(p-q)a + 2pqd$	
additive	$2q\alpha$	$(q-p)\alpha$	$-2p\alpha$	$\alpha=a+(q-p)d$
dominance dev.	$-q^2d$	$2pqd$	$-p^2d$	

$$\text{Mean(additive effect)} = 0, \text{mean(dominance deviation)} = 0, \text{cov(additive effect, dominance dev)} = 0$$

$$\begin{aligned} \text{Genetic variance} = V_G &= 2pq\alpha^2 + (2pqd)^2 \\ &= V_A + V_D \end{aligned}$$

22

## Covariance between genetic effects of relatives

Model with 1 gene, 2 alleles and additive and dominant gene action

$$\text{Cov}(g_i, g_j) = \text{Cov}(a_i + d_i, a_j + d_j) = \text{Cov}(a_i, a_j) + \text{cov}(d_i, d_j) \\ = A(i,j) V_A + D(i,j) V_D$$

$D(i,j)$  = prob(i and j inherit the same genotype IBD)

Eg

$D(i,j) = 1$  for MZ twins, 0.25 for full-sibs, 0 for parent and offspring

23

## Covariance between genetic effects of relatives

Model with 1 gene, 2 alleles and additive and dominant gene action

Relationships	MZ twins	full-sibs	1/2sibs	P-O
A	1	0.5	0.25	0.5
D	1	0.25	0	0

Therefore can estimate both  $V_A$  and  $V_D$  by using multiple relationships

24

## Covariance between environmental effects of relatives

$y = \text{mean} + \text{genetic effect} + \text{common environment effect} + \text{individual environment effect}$

$$y = \text{mean} + g + e_c + e$$

Model with a common environmental effect within the same family

$\text{Cov}(e_{ci}, e_{cj}) = V_c$  if i and j in same family, zero otherwise

Relationships	MZ twins	full-sibs	1/2sibs	P-O
A	1	0.5	0.25	0.5
D	1	0.25	0	0
E common	1	1	?	?

25

## Covariance between relatives

Estimating  $V_A$ ,  $V_D$  and  $V_C$

Difficult!

Assume  $V_D = 0$

$$VA = 2(\text{cov(MZ twins)} - \text{cov(full-sibs)})$$

Relationships	MZ twins	full-sibs	1/2sibs	P-O
A	1	0.5	0.25	0.5
D	1	0.25	0	0
E common	1	1	?	?

26

# Covariance between relatives

Can add epistatic interactions to model

$$g = \text{mean} + \text{additive} + \text{dominance} + \text{epistasis}$$

$$\text{eg } g = \text{mean} + a + d + aa$$

$$\text{Cov}(g_i, g_j) = A(i,j) V_A + D(i,j) V_D + A(i,j)^2 V_{AA}$$

Relationships	MZ twins	full-sibs	1/2sibs	P-O
A	1	0.5	0.25	0.5
D	1	0.25	0	0
AxA	1	0.25	0.0625	0.25

27

VOLUME II

NOVEMBER, 1903

No. 4

## ON THE LAWS OF INHERITANCE IN MAN\*.

### I. INHERITANCE OF PHYSICAL CHARACTERS.

By KARL PEARSON, F.R.S., assisted by ALICE LEE, D.Sc.  
University College, London.

364      *On the Laws of Inheritance in Man*

DIAGRAM IV. Distribution of Stature.

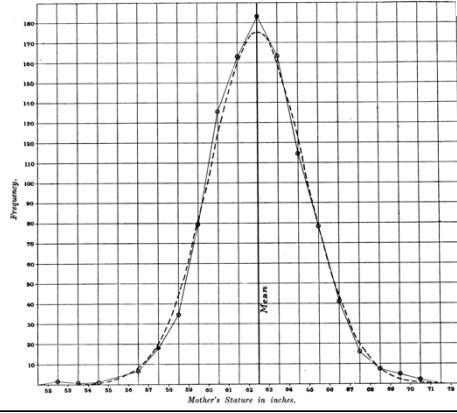
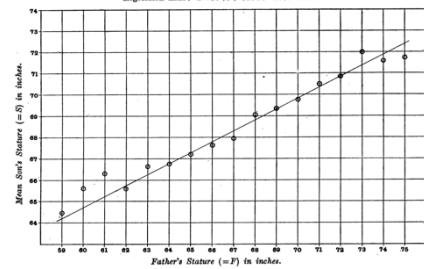


DIAGRAM I. Probable Stature of Son for given Father's Stature.

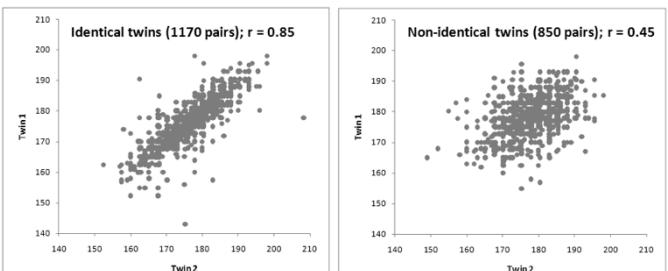
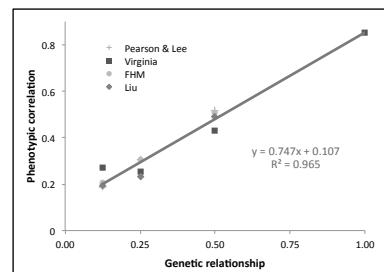
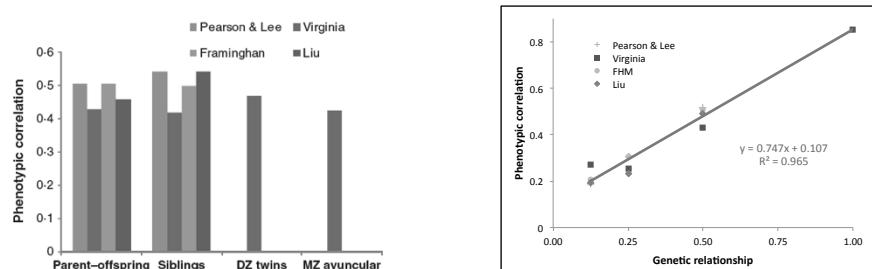
Regression Line:  $S = 53.73 + .516 F$ . 1078 Cases.



PAIR	CORRELATION	SE
Spouse	0.28	0.02
Son-Father	0.51	0.02
Daughter-Father	0.51	0.01
Son-Mother	0.49	0.02
Daughter-Mother	0.51	0.01
Brother-brother	0.51	0.03
Sister-sister	0.54	0.02
Brother-sister	0.55	0.01

28

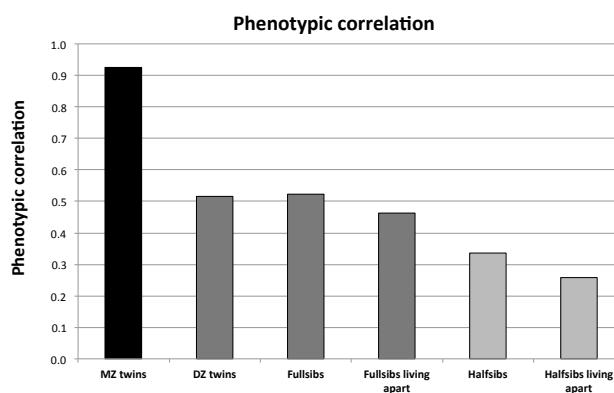
## Resemblance between relatives (height)



29

## More data on height

Data from ~172,000 18-year old brother pairs



[Magnus Johannesson, David Cesarini] 30

# Sex Differences in Heritability of BMI: A Comparative Study of Results from Twin Studies in Eight Countries

Twin Research October 2003

Karoline Schousboe<sup>1</sup>, Gonke Willemsen<sup>2</sup>, Kirsten O. Kyvik<sup>1</sup>, Jakob Mortensen<sup>1</sup>, Dorret I. Boomsma<sup>2</sup>, Belinda K. Cornes<sup>3</sup>, Chayna J. Davis<sup>4</sup>, Corrado Fagnani<sup>5</sup>, Jacob Hjelmborg<sup>6</sup>, Jaakko Kaprio<sup>6</sup>, Marlies de Lange<sup>7</sup>, Michelle Luciano<sup>8</sup>, Nicholas G. Martin<sup>9</sup>, Nancy Pedersen<sup>10</sup>, Kirsi H. Pietiläinen<sup>4,8</sup>, Aila Rissanen<sup>8</sup>, Suoma Saarni<sup>8</sup>, Thorkild I.A. Sørensen<sup>9</sup>, G. Caroline M. van Baal<sup>2</sup>, and Jennifer R. Harris<sup>10</sup>

**Table 5a**

Twin Correlations (R) for BMI and Number of Pairs (N) Assessed by Zygosity and Sex for Twins Aged 20–29 years

	Australia R (N)	Denmark R (N)	Finland R (N)	Italy R (N)	Netherlands R (N)	Norway R (N)	Sweden R (N)	UK R (N)
MZm	0.67 (390)	0.77 (824)	0.74 (247)	0.83 (66)	0.65 (299)	0.69 (563)	0.77 (887)	n.a.
DZm	0.32 (260)	0.35 (897)	0.32 (304)	0.52 (43)	0.31 (222)	0.41 (479)	0.35 (1346)	n.a.
MZf	0.72 (768)	0.73 (1161)	0.78 (411)	0.83 (129)	0.79 (518)	0.74 (738)	0.73 (1054)	0.74 (89)
DZf	0.33 (496)	0.35 (1046)	0.37 (358)	0.58 (76)	0.41 (336)	0.35 (643)	0.36 (1472)	0.52 (75)
DZOS	0.18 (596)	0.30 (1620)	0.22 (668)	0.12 (96)	0.36 (473)	0.18 (968)	n.a.	n.a.

## Average correlations

MZ	0.74
DZ (same sex)	0.36
DZ (opposite sex)	0.25

31

# Variability in the heritability of body mass index: a systematic review and meta-regression

Cathy E. Elks<sup>1</sup>, Marcel den Hoed<sup>1</sup>, Jing Hua Zhao<sup>1</sup>, Stephen J. Sharp<sup>1</sup>, Nicholas J. Wareham<sup>1</sup>, Ruth J. F. Loos<sup>1</sup> and Ken K. Ong<sup>1,2\*</sup>

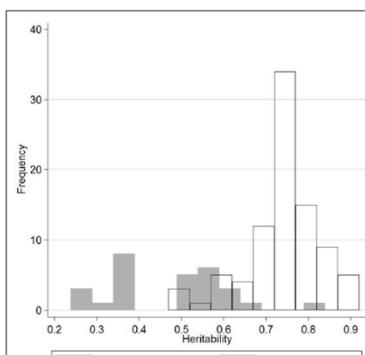
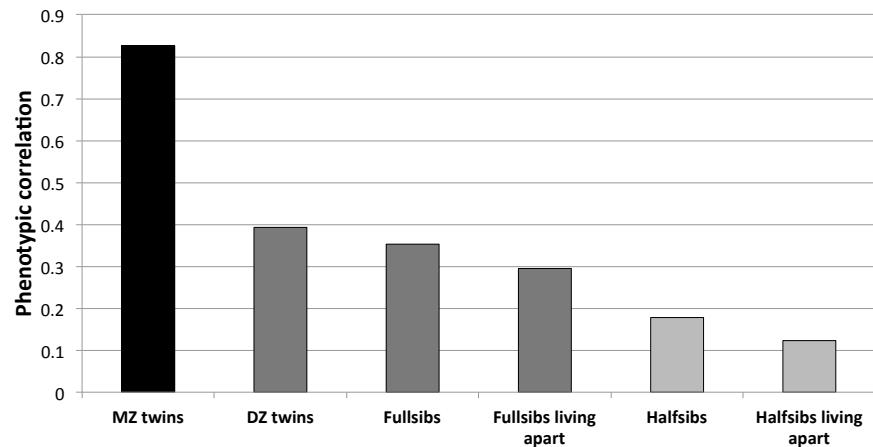


FIGURE 1 | Histogram showing the wide distribution of reported estimates of BMI heritability from twin studies (white bars) and family studies (gray bars).

32

# BMI

Data from ~172,000 18-year old brother pairs



[Magnus Johannesson, David Cesari] 33

# Estimating genetic variation within families

Peter M. Visscher  
[peter.visscher@uq.edu.au](mailto:peter.visscher@uq.edu.au)

1

## Key concepts

1. There is variation in realised relationships given the expected value from the pedigree;
2. Variation in realised relationships can be captured with genetic markers;
3. Variation in realised relationships can be exploited to estimate genetic variation

2

## Genetic covariance between relatives

$$\text{cov}_G(y_i, y_j) = a_{ij}\sigma_A^2 + d_{ij}\sigma_D^2$$

$a =$  additive coefficient of relationship  
 $= 2 * \theta (= E(\pi_a))$

$d =$  coefficient of fraternity  
 $= \text{Prob}(2 \text{ alleles are IBD}) = \Delta = E(\pi_d)$

3

## Examples (no inbreeding)

Relatives	a	d
MZ twins	1	1
Parent-offspring	$\frac{1}{2}$	0
Fullsibs	$\frac{1}{2}$	$\frac{1}{4}$
Double first cousins	$\frac{1}{4}$	$\frac{1}{16}$

4

## Controversy/confounding: nature vs nurture

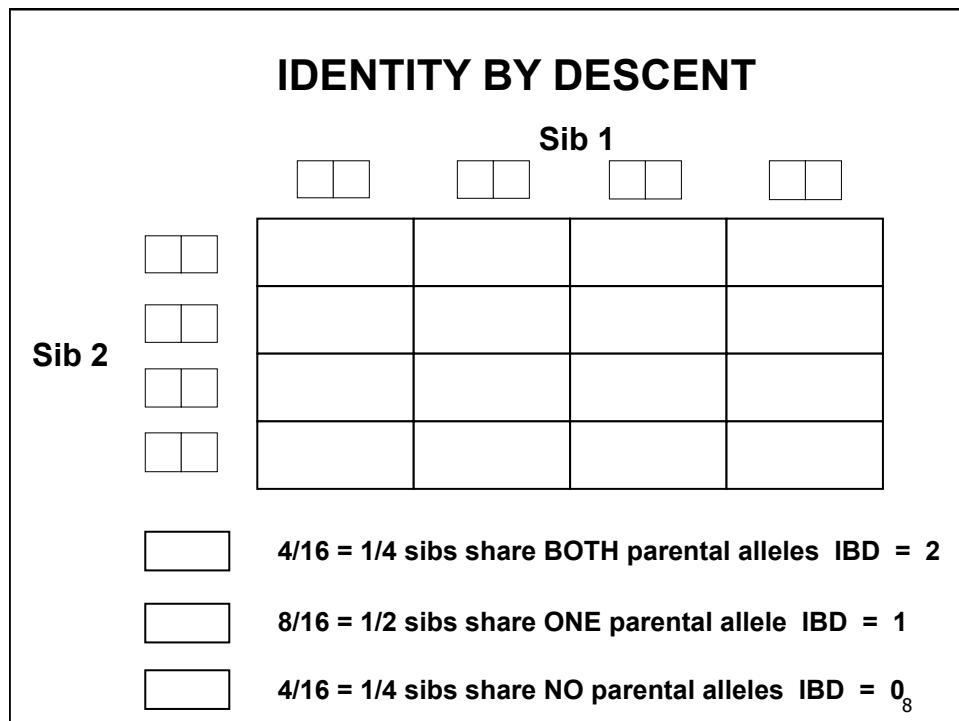
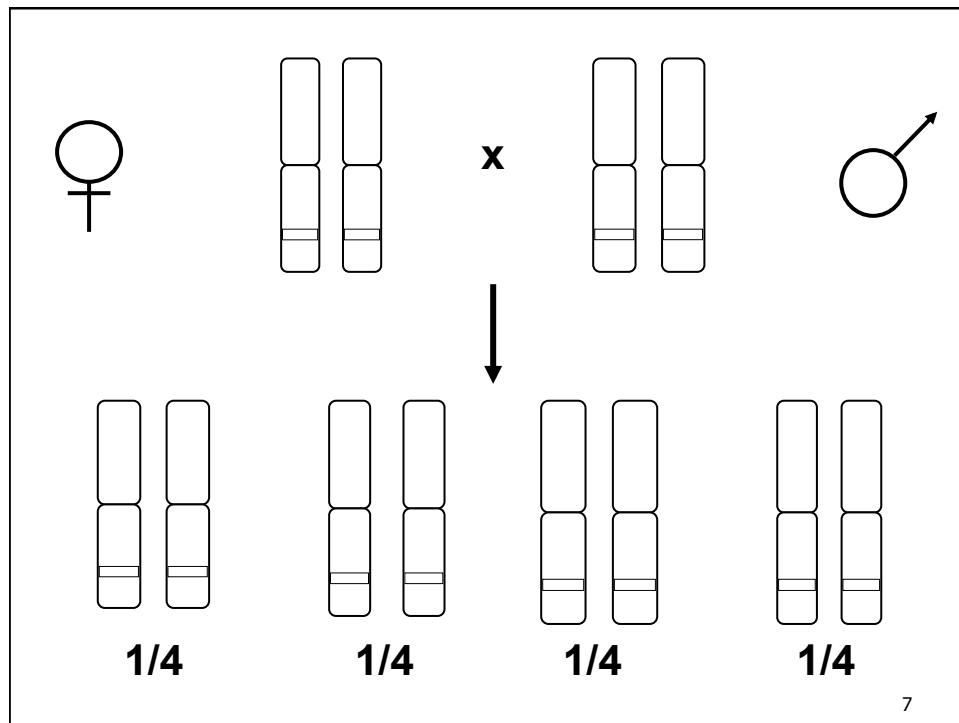
- Is observed resemblance between relatives genetic or environmental?
  - MZ & DZ twins (shared environment)
  - Fullsibs (dominance & shared environment)
- Estimation and statistical inference
  - Different models with many parameters may fit data equally well



## *Actual or realised genetic relationship*

= proportion of genome shared IBD ( $\pi_a$ )

- Varies around the expectation
  - Apart from parent-offspring and MZ twins
- Can be estimated using marker data



## Single locus

Relatives	$E(\pi_a)$	$\text{var}(\pi_a)$
Fullsibs	$\frac{1}{2}$	$\frac{1}{8}$
Halfsibs	$\frac{1}{4}$	$\frac{1}{16}$
Double 1 <sup>st</sup> cousins	$\frac{1}{4}$	$\frac{3}{32}$

9

## Several notations

IBD	Probability	Actual	Realisations
IBD0	$k_0$	0 or 1	$k_0 \quad k_1 \quad k_2$
IBD1	$k_1$	0 or 1	$1 \quad 0 \quad 0$
IBD2	$k_2$	0 or 1	$0 \quad 1 \quad 0$
	$\Sigma=1$	$\Sigma=1$	$0 \quad 0 \quad 1$

$$\pi_a = \frac{1}{2}k_1 + k_2 = R = 2\theta$$

$$\pi_d = k_2 = \Delta_{xy}$$

10

[e.g., LW Chapter 7; Weir and Hill 2011, Genetics Research]

## $n$ multiple unlinked loci

Relatives	$E(\pi_a)$	$\text{var}(\pi_a)$
Fullsibs	$\frac{1}{2}$	$\frac{1}{8n}$
Halfsibs	$\frac{1}{4}$	$\frac{1}{16n}$
Double 1 <sup>st</sup> cousins	$\frac{1}{4}$	$\frac{3}{32n}$

11

## Loci are on chromosomes

- Segregation of large chromosome segments within families
  - increasing variance of IBD sharing
- Independent segregation of chromosomes
  - decreasing variance of IBD sharing

12

## Theoretical SD of $\pi_a$

Relatives	1 chrom (1 M)	genome (35 M)
FULLSIBS	0.217	0.038
HALFSIBS	0.154	0.027
Double 1 <sup>st</sup> cousins	0.173	0.030

[Stam 1980; Hill 1993; Guo 1996; Hill & Weir 2011]<sup>13</sup>

## FULLSIBS: genome-wide (Total length L Morgan)

$$\text{var}(\pi_a) \approx 1/(16L) - 1/(3L^2) \quad [\text{Stam 1980; Hill 1993; Guo 1996}]$$

$$\text{var}(\pi_d) \approx 5/(64L) - 1/(3L^2)$$

$$\text{var}(\pi_d)/\text{var}(\pi_a) \approx 1.3 \text{ if } L = 35$$

Genome-wide variance depends more on total genome length than on the number of chromosomes

## Fullsibs: Correlation additive and dominance relationships

$$r(\pi_a, \pi_d) = \sigma(\pi_a) / \sigma(\pi_d) \approx [1/(16L) / (5/(64L))]^{0.5} = 0.89.$$

Using  $\beta(\pi_a \text{ on } \pi_d) = 1$

Difficult but not impossible to disentangle additive and dominance variance

NB Practical

15

## Summary Additive and dominance (fullsibs)

	SD( $\pi_a$ )	SD( $\pi_d$ )
Single locus	0.354	0.433
One chromosome (1M)	0.217	0.247
Whole genome (35M)	0.038	0.043
Predicted correlation (genome-wide $\pi_a$ and $\pi_d$ )		0.89

16

## Estimating IBD from marker data

- Elston-Stewart algorithm  
Handles large pedigrees, but small nr of loci, exact IBD distributions (Elston and Stewart, 1971)
- Lander-Green algorithm  
Handles small pedigrees, but large nr of loci, exact IBD distributions (Lander and Green, 1987). Software: Merlin
- MCMC methods  
Calculates approximate IBD distributions (Heath, 1997). Software: Loki
- Average sharing methods.  
Calculates approximate IBD distributions (Fulker et al., 1995; Almasy and Blangero, 1998). Software: SOLAR

17

## Estimating $\pi$ when marker is not fully informative

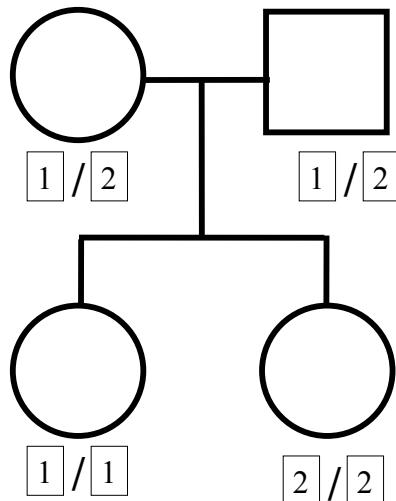
- Using:
  - Mendelian segregation rules
  - Marker allele frequencies in the population

18

IBD can be trivial...

IBD=0

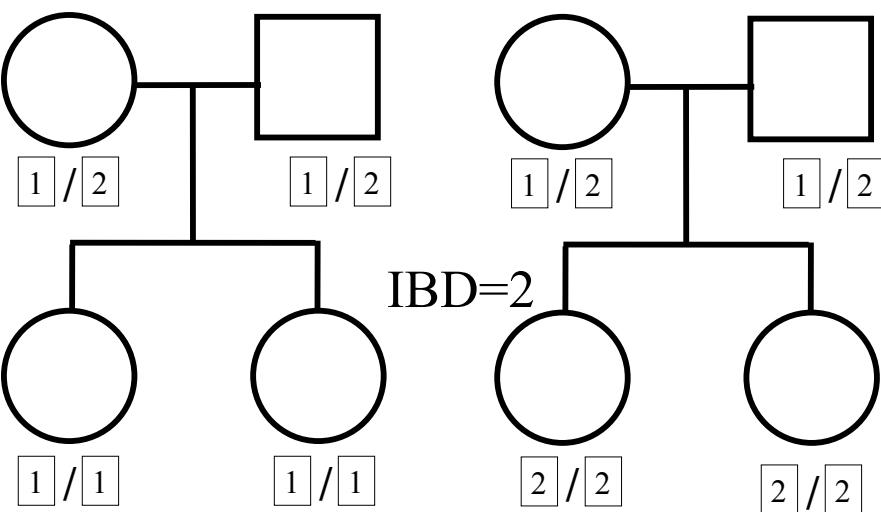
19



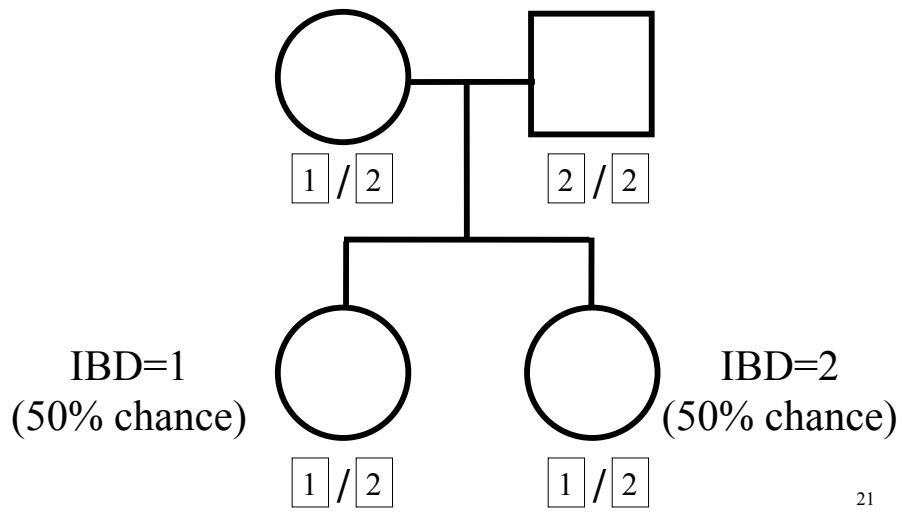
Two Other Simple Cases...

IBD=2

20

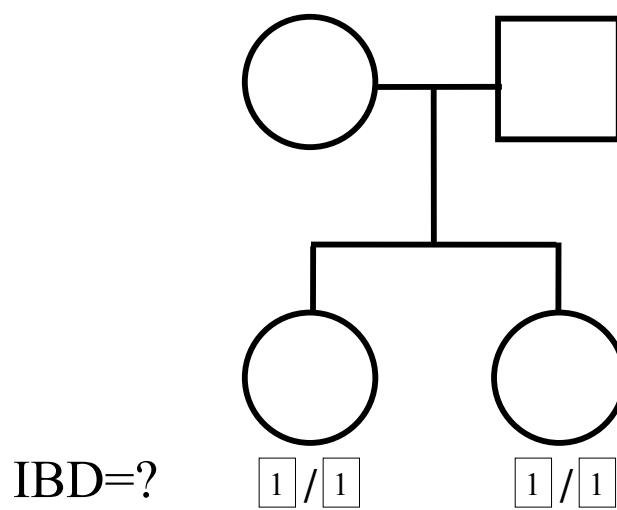


A little more complicated...



21

And even more complicated...



22

## Bayes Theorem for IBD Probabilities

$$\begin{aligned}
 P(\text{IBD} = i | G) &= \frac{P(\text{IBD} = i, G)}{P(G)} \\
 \text{prior} \longrightarrow &= \frac{P(\text{IBD} = i)P(G | \text{IBD} = i)}{P(G)} \\
 &= \frac{P(\text{IBD} = i)P(G | \text{IBD} = i)}{\sum_j P(\text{IBD} = j)P(G | \text{IBD} = j)}
 \end{aligned}$$

posterior

Prob(data)

23

## P(Marker Genotype|IBD State)

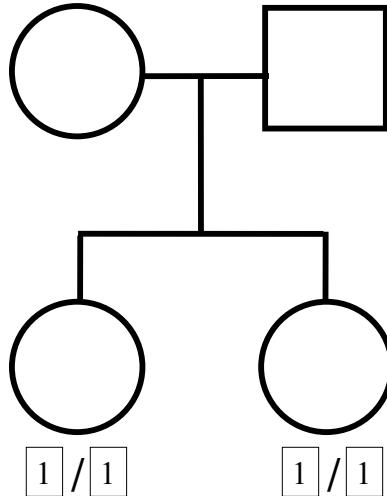
Sib	CoSib	IBD		
		0	1	2
(a,b)	(c,d)	$p_a p_b p_c p_d$	0	0
(a,a)	(b,c)	$p_a^2 p_b p_c$	0	0
(a,a)	(b,b)	$p_a^2 p_b^2$	0	0
(a,b)	(a,c)	$p_a^2 p_b p_c$	$p_a p_b p_c$	0
(a,a)	(a,b)	$p_a^3 p_b$	$p_a^2 p_b$	0
(a,b)	(a,b)	$p_a^2 p_b^2$	$p_a p_b^2 + p_a^2 p_b$	$p_a p_b$
(a,a)	(a,a)	$p_a^4$	$p_a^3$	$p_a^2$
Prior Probability		$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

[Assumes Hardy-Weinberg proportions of genotypes in the population]

24

$$p_1 = 0.5$$

## Worked Example



$$P(G | IBD = 0) = p_1^4 = \frac{1}{16}$$

$$P(G | IBD = 1) = p_1^3 = \frac{1}{8}$$

$$P(G | IBD = 2) = p_1^2 = \frac{1}{4}$$

$$P(G) = \frac{1}{4} p_1^4 + \frac{1}{2} p_1^3 + \frac{1}{4} p_1^2 = \frac{9}{64}$$

$$P(IBD = 0 | G) = \frac{\frac{1}{4} p_1^4}{P(G)} = \frac{1}{9}$$

$$P(IBD = 1 | G) = \frac{\frac{1}{2} p_1^3}{P(G)} = \frac{4}{9}$$

$$\hat{\pi} = \frac{2}{3}$$

$$P(IBD = 2 | G) = \frac{\frac{1}{4} p_1^2}{P(G)} = \frac{4}{9}$$

25

## Application (1)

**Aim: estimate genetic variance from actual relationships between fullsib pairs**

- Two cohorts of Australian twin families

	Adolescent	Adult
Families	500	1512
Individuals	1201	3804
Sibpairs with genotypes	950	3451
Markers per individual	211-791	201-1717
Average marker spacing	6 cM	5 cM

26

## Application (1)

- Phenotype = height

**Number of sibpairs with phenotypes and genotypes**

<i>Adolescent cohort</i>	931
<i>Adult cohort</i>	2444
<i>Combined</i>	3375

27

**Mean IBD sharing across the genome for the  $j$ th sib pair was based on IBD estimated every centimorgan and averaged over 3500 points ( $L = 35$ )**

$$\text{additive} \quad \bar{\hat{\pi}}_{a(j)} = \sum_{i=1}^{3500} \hat{\pi}_{a(ij)} / 3500$$

$$\text{dominance} \quad \bar{\hat{\pi}}_{d(j)} = \sum_{i=1}^{3500} p_{2(ij)} / 3500$$

28

**And for the  $c^{\text{th}}$  chromosome of length  $l_c$  cM**

additive

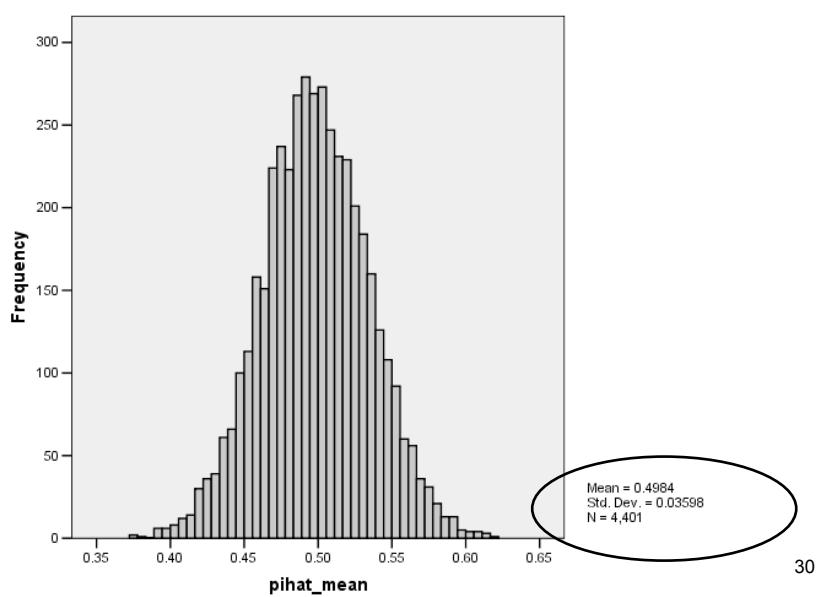
$$\bar{\hat{\pi}}_{a(j)}^c = \sum_{i=1}^{l_c} \hat{\pi}_{a(ij)}^c / l_c$$

dominance

$$\bar{\hat{\pi}}_{d(j)}^c = \sum_{i=1}^{l_c} p_{2(ij)} / l_c$$

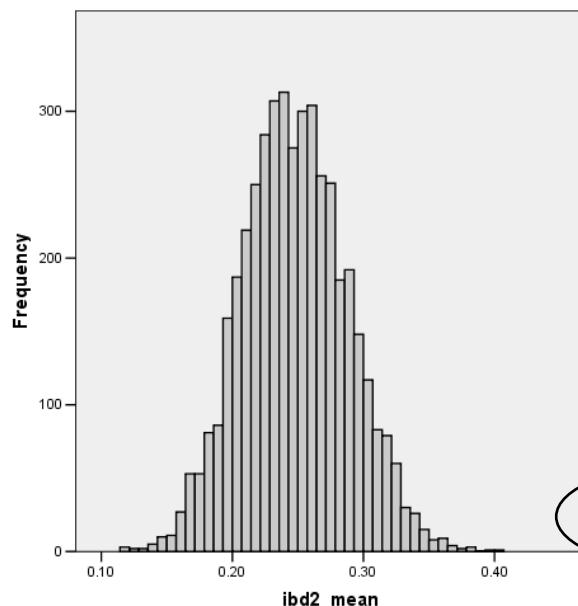
29

Mean and SD of genome-wide additive relationships



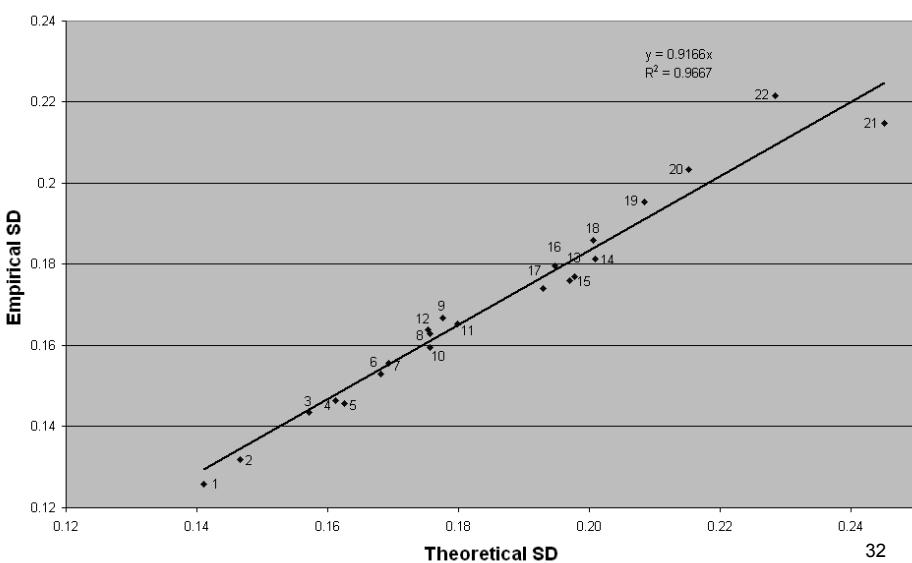
30

### Mean and SD of genome-wide dominance relationships



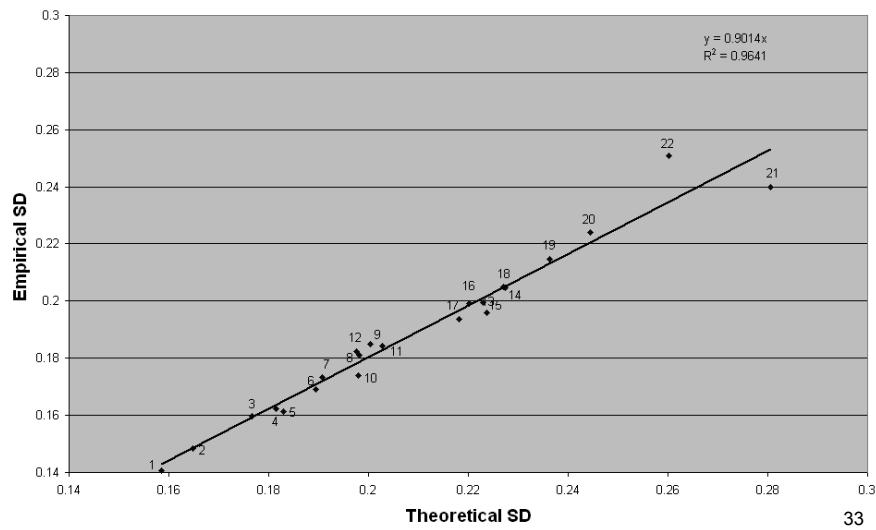
31

### Empirical and theoretical SD of additive relationships correlation = 0.98 ( $n = 4401$ )



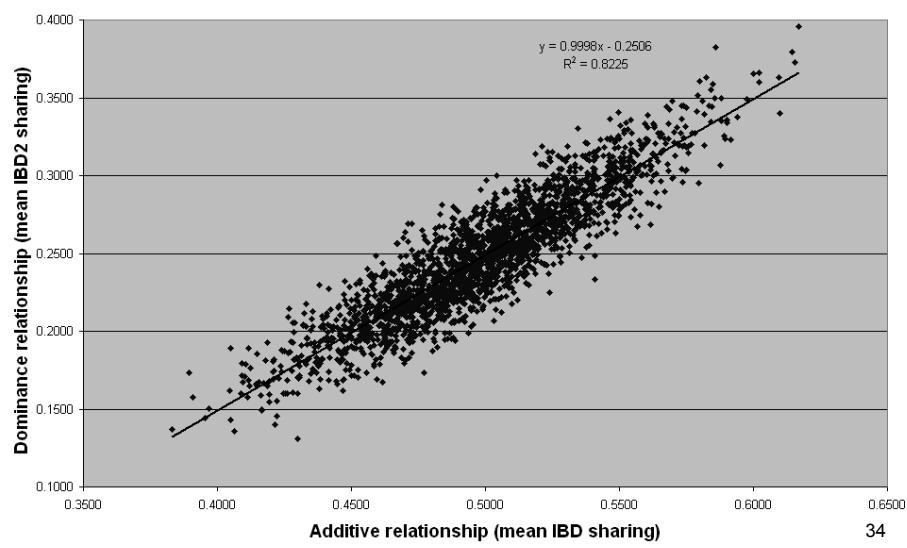
32

Empirical and theoretical SD of dominance relationships  
correlation = 0.98 ( $n = 4401$ )



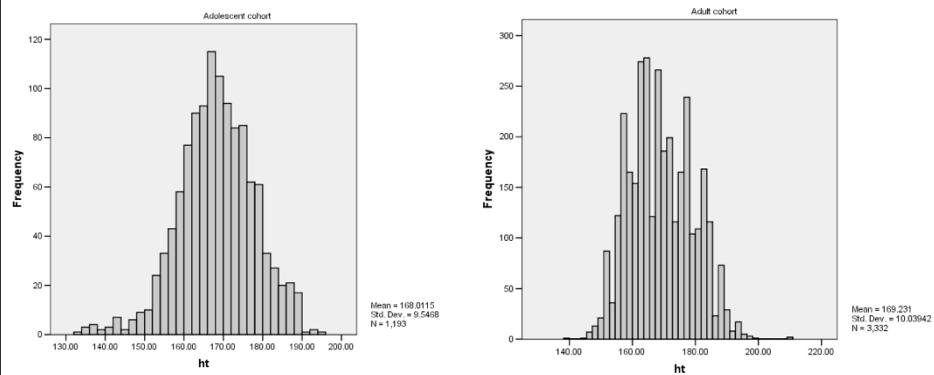
33

Additive and dominance relationships  
correlation = 0.91 ( $n= 4401$ )



34

# Phenotypes



After adjustment for sex and age:

$$\sigma_p = 7.7 \text{ cm}$$

$$\sigma_p = 6.9 \text{ cm}$$

35

## Phenotypic correlation between siblings

	Raw	After age & sex
Adolescents	0.33	0.40
Adults	0.24	0.39

36

## Models

$$y_{ij} = \mu + c_i + a_{ij} + e_{ij}$$

$$\text{var}(y) = \sigma_c^2 + \sigma_a^2 + \sigma_e^2$$

$$\text{cov}(y_{ij}, y_{ik}) = \sigma_c^2 + \pi_{a(jk)} \sigma_a^2$$

C = Family effect

A = Genome-wide additive genetic

E = Residual

Full model              C + A + E

Reduced model          C + E

37

## Estimation

- Maximum Likelihood variance components
- Likelihood-ratio-test (LRT) to calculate P-values for hypotheses

$$H_0: A = 0$$

$$H_1: A > 0$$

38

## Estimates: null model (CE)

Cohort	Family effect (C)
<i>Adolescent</i>	0.40 (0.34 – 0.45)
<i>Adult</i>	0.39 (0.36 – 0.43)
<i>Combined</i>	0.39 (0.36 – 0.42)

39

## Estimates: full model (ACE)

Cohort	C	A	P
<i>Adolescent</i>	0	0.80	0.0869
<i>Adult</i>	0	0.80	0.0009
<i>Combined</i>	0	0.80	0.0003

► **All family resemblance due to additive genetic variation**

40

## Sampling variances are large

Cohort	A (95% CI)
<i>Adolescent</i>	0.80 (0.00 – 0.90)
<i>Adult</i>	0.80 (0.43 – 0.86)
<i>Combined</i>	0.80 (0.46 – 0.85)

41

## Power and SE of estimates

- True parameter ( $t$  = intra-class correlation)
- Sample size ( $n$  pairs)
- Variance in genome-wide IBD sharing ( $\text{var}(\pi)$ )

$$\text{var}(\hat{h}^2) \approx (1 - t^2)^2 / [(1 + t^2)(n \text{var}(\pi))]$$

$$NCP = nh^4 \text{var}(\pi)(1+t^2) / (1-t^2)^2$$

42

## Application (2)

### Genome partitioning of additive genetic variance for height

- Aims
  - Estimate genetic variance from genome-wide IBD in larger sample
  - Partition genetic variance to individual chromosomes
    - using chromosome-wide coefficients of relationship
  - Test hypotheses about the distribution of genetic variance in the genome

43

<b>Sample</b>	<b># Sibpairs</b>	<b>Sib Correlation</b>
AU	5952	0.43
US	3996	0.50
NL	1266	0.45
Total	11,214	0.46

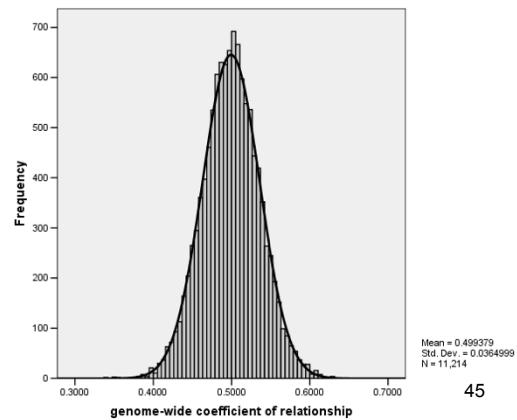
44

# Realised relationships

Mean 0.499

Range 0.31 – 0.64

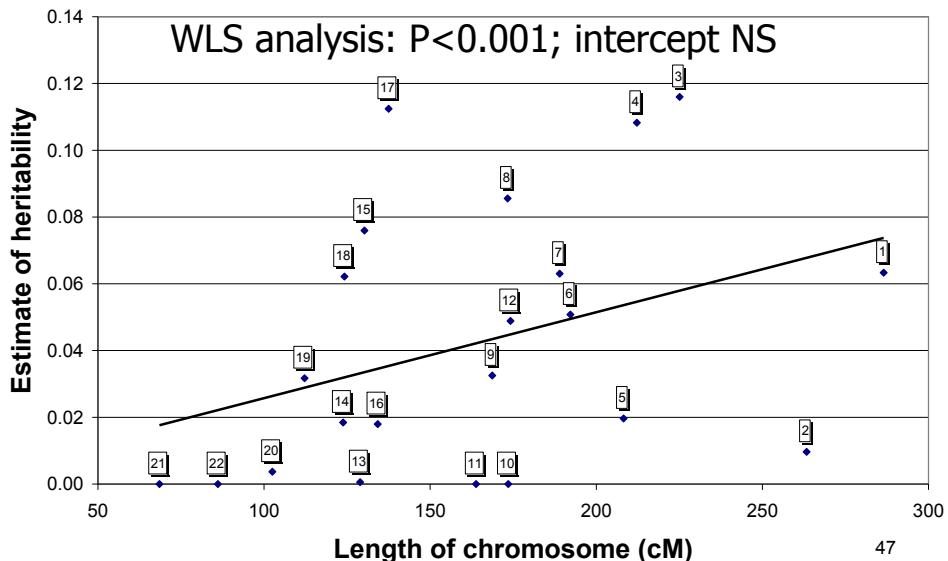
SD 0.036



45

Chrom.	Single chromosome analyses					Combined chromosome analysis		
	$f^2$ (a)	$h_i^2$ (b)	$e^2$ (c)	LRT <sup>d</sup>	P-value	$h_i^2$	LRT <sup>e</sup>	P-value
1	0.4285	0.0607	0.5108	1.201	0.137	0.0633	1.418	0.117
2	0.4525	0.0131	0.5344	0.065	0.399	0.0097	0.037	0.424
3	0.4023	0.1134	0.4843	5.704	0.008	0.1160	6.269	0.006
4	0.4036	0.1124	0.4840	5.938	0.007	0.1082	5.705	0.008
5	0.4458	0.0264	0.5278	0.319	0.286	0.0196	0.191	0.500
6	0.4336	0.0506	0.5158	1.294	0.128	0.0508	1.370	0.500
7	0.4284	0.0616	0.5100	2.019	0.078	0.0630	2.230	0.068
8	0.4234	0.0708	0.5058	2.778	0.048	0.0856	4.172	0.021
9	0.4482	0.0216	0.5302	0.277	0.299	0.0325	0.663	0.500
10	0.4590	0.0000	0.5410	0.000	0.500	0.0000	0.000	0.500
11	0.4590	0.0000	0.5410	0.000	0.500	0.0000	0.000	0.500
12	0.4365	0.0451	0.5184	1.121	0.145	0.0489	1.434	0.500
13	0.4545	0.0089	0.5366	0.056	0.406	0.0006	0.000	0.500
14	0.4427	0.0323	0.5250	0.728	0.197	0.0185	0.246	0.500
15	0.4241	0.0703	0.5056	3.353	0.034	0.0760	4.028	0.022
16	0.4556	0.0069	0.5375	0.035	0.426	0.0180	0.251	0.308
17	0.4023	0.1142	0.4834	9.019	0.001	0.1124	8.967	0.001
18	0.4237	0.0703	0.5060	3.753	0.026	0.0622	3.013	0.041
19	0.4437	0.0309	0.5253	0.759	0.192	0.0317	0.840	0.500
20	0.4575	0.0031	0.5395	0.008	0.464	0.0037	0.012	0.456
21	0.4590	0.0000	0.5410	0.000	0.500	0.0000	0.000	0.500
22	0.4590	0.0000	0.5410	0.000	0.500	0.0000	0.000	0.500
SUM		0.9126		38.427		0.9205	40.846	46

## Longer chromosomes explain more additive genetic variance: ~0.03 per 100 cM



## Application (3)

- Using SNP data to estimate IBD
- Data from ~20,000 fullsib pairs
- Height and BMI

**frontiers in  
ENDOCRINOLOGY**

ORIGINAL RESEARCH ARTICLE  
published: 28 February 2012  
doi: 10.3389/fendo.2012.00029

Variability in the heritability of body mass index: a systematic review and meta-regression

Cathy E. Elks<sup>1</sup>, Marcel den Hoed<sup>1</sup>, Jing Hua Zhao<sup>1</sup>, Stephen J. Sharp<sup>1</sup>, Nicholas J. Wareham<sup>1</sup>, Ruth J. F. Loos<sup>1</sup> and Ken K. Ong<sup>1,2\*</sup>

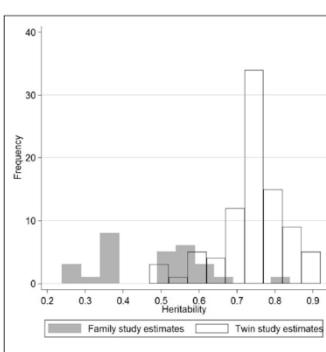
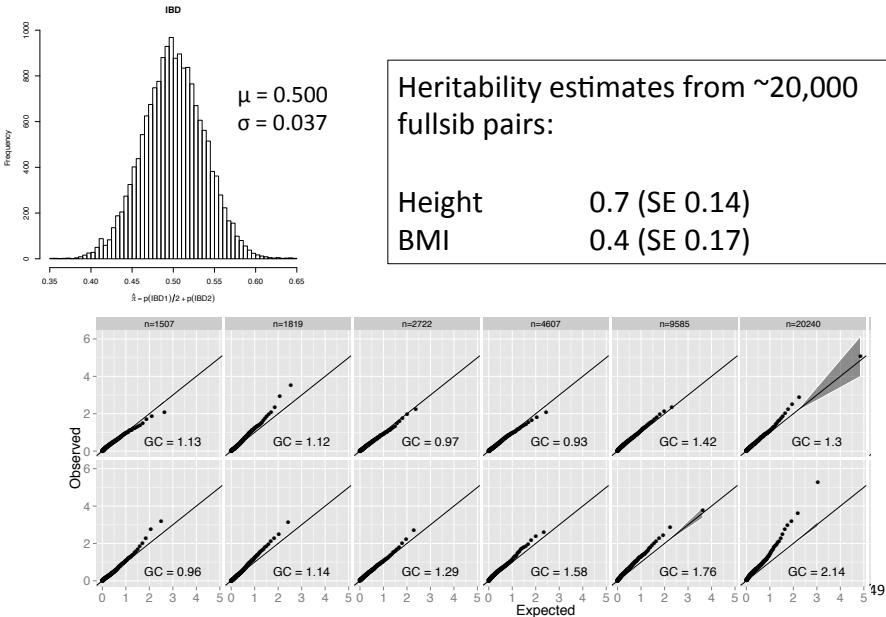


FIGURE 1 | Histogram showing the wide distribution of reported estimates of BMI heritability from twin studies (white bars) and family studies (gray bars).

## Genetic variation within families using SNP data



## Conclusions

- Empirical variation in genome-wide IBD sharing follows theoretical predictions
- Genetic variance can be estimated from genome-wide IBD within families
  - results for height consistent with estimates from between-relative comparisons
  - no assumptions about nature/nurture causes of family resemblance
- Genetic variance can be partitioned onto chromosomes

## Key concepts

1. There is variation in realised relationships given the expected value from the pedigree;
2. Variation in realised relationships can be captured with genetic markers;
3. Variation in realised relationships can be exploited to estimate genetic variation

# Estimating relationship from marker genotypes

Mike Goddard

1

## Relationships

We use relationship data  
to estimate genetic variance  
to estimate demographic history  
...

2

## Relationships

Additive genetic relationship  $G(i, j)$

= proportion of the genome in i and j that  
is IBD

Pedigree relationship  $A(i,j) = \text{Prob}(\text{IBD})$

=  $E(G(i,j))$

Actual relationship deviates randomly from this  
expectation

3

## Relationships

Single locus case, full sibs

Parents  $A_1A_2 \times A_3A_4$

offspring       $A_1A_3$   
                   $A_1A_4$   
                   $A_2A_3$   
                   $A_2A_4$

Pairs of sibs share

0 alleles	25% of the time
1 allele	50%
2 alleles	25%

$E(G) = A = 0.5$  but  $G$  varies from 0 to 1

4

## Estimate relationship from markers

G is a more accurate description of relationship than A

- G captures unknown pedigree information
- pedigree can be incorrect
- G captures deviations from A

Therefore, can use G in

- Random sample of population ("unrelated individuals")
- Individuals with same pedigree

5

## Estimate relationship from markers

1. Well defined (recent) base
2. No well defined base
3. Well defined, recent base

Eg Data on families of full-sibs and parents of sibs are the base

6

## Estimate relationship from markers

Eg Data on families of full-sibs and parents of sibs are the base

Consider a single SNP

Full sibs can be IBD at either maternal or paternal allele

	IBD status	P(IBD status)
Maternal	Paternal	
yes	yes	0.25
yes	no	0.25
no	yes	0.25
no	no	0.25

7

## Estimate relationship from markers

Eg Data on families of full-sibs and parents of sibs are the base

At this SNP, one sib has genotype AA and the other is AB, mother = AB, father = AA

$$\begin{aligned} P(\text{IBD status} \mid \text{SNP genotypes}) &= \frac{P(\text{SNP genotypes} \mid \text{IBD status}) * P(\text{IBD status})}{P(\text{SNP genotypes})} \\ &= \frac{P(\text{SNP genotypes} \mid \text{IBD status}) * P(\text{IBD status})}{\sum P(\text{SNP genotypes} \mid \text{IBD status}) * P(\text{IBD status})} \end{aligned}$$

8

## Estimate relationship from markers

$$= \frac{P(\text{SNP genotypes} | \text{IBD status}) * P(\text{IBD status})}{\sum P(\text{SNP genotypes} | \text{IBD status}) * P(\text{IBD status})}$$

IBD status		P(IBD status)	P(genotypes IBD status)	P(IBD status   genotypes)	G
Maternal	Paternal				
yes	yes	0.25	0	0	1
yes	no	0.25	0	0	0.5
no	yes	0.25	1	0.5	0.5
no	no	0.25	1	0.5	0

$$\sum P(\text{SNP genotypes} | \text{IBD status}) * P(\text{IBD status}) = 0.5$$

$$E(G) = 0.25 \text{ compared with } A=0.5$$

9

## Estimate relationship from markers

1. Well defined, recent base

Eg Data on families of full-sibs and parents of sibs are the base

- a) Calculate Bayesian probability of IBD status at each SNP  
→ E(G) at each SNP  
average over SNPs

- b) Use haplotypes ?

10

## Estimate relationship from markers

2. Less well defined, less recent base

Eg Data on current population, base = ancestors 1000 years ago  
and allele frequencies in base are known ( $p$  and  $q$ )

Consider haploid gametes of SNP alleles instead of genotypes  
What fraction of the gametes are IBD ( $G$ )?

At a single SNP, there are 3 possible data sets and their probabilities are

A and A	A and B	B and B
$p^2 + pqG$	$2pq(1-G)$	$q^2 + pqG$

11

## Estimate relationship from markers

SNP genotypes	A and A	A and B	B and B
Probability	$p^2 + pqG$	$2pq(1-G)$	$q^2 + pqG$
score ( $x$ )	$q/p$	-1	$p/q$

Estimate  $G(i,j)$  from the mean value of  $x$  over SNPs

This is a relationship between gametes. Calculate  $G$  for individuals from the 4 gametic relationships.

See Yang et al (2010) and Powell et al (2010) for the diploid formulae.

12

## Estimate relationship from markers

### 2. No well defined base

Eg random sample from population but don't know allele frequency in the base.

a) Use the current population as the base

Problem: Some  $G < 0$

Cannot interpret as probabilities but still interpret as covariances

If  $g$  = genetic value,  $V(g) = G V_A$

where  $G$  is calculated as above but using allele frequencies in current population.

13

## Estimate relationship from markers

### 2. No well defined base

b) Assume SNPs are a random sample of loci as are QTL

$$y = \text{mean} + g + e$$

$$y = \text{mean} + Zu + e$$

$Z_{ij} = 0$  for AA, 1 for AB or 2 for BB

$u \sim N(0, I\sigma_u^2) \rightarrow g = Zu \sim N(0, ZZ'\sigma_u^2)$ ,  $ZZ'\sigma_u^2 = G\sigma_g^2$ , if  $\sigma_g^2 = N\sigma_u^2$

where  $N$ =number of SNPs

Therefore,  $G = ZZ'/N$

14

## Estimate relationship from markers

2a and 2b compared for gametic relationships

SNP data	A and A	A and B	B and B
score (x)	$q/p$	-1	$p/q$
weight (w)	$pq$	$pq$	$pq$

2a)  $G = \text{mean of } x$

2b)  $G = \text{weighted mean of } x = \sum wx / \sum w$

This could be described as using the IBS status of SNPs instead of IBD

15

## Estimate relationship from markers

2a)  $G = \text{mean of } x$

gives more emphasis to sharing rare alleles

Makes sense because individuals who share rare alleles are more likely to be closely related than individuals who share common alleles.

Gives minimum error variance of relationship under some conditions

16

## Estimate relationship from markers

2. No well defined base

c) Assume SNPs are a random sample of loci as are QTL but effect of SNP decreases as heterozygosity increases

$$y = \text{mean} + g + e$$

$$y = \text{mean} + Zu + e$$

$Z_{ij} = 0$  for AA, 1 for AB or 2 for BB

$u \sim N(0, D\sigma_u^2) \rightarrow g = Zu \sim N(0, ZDZ'\sigma_u^2)$ ,  $ZDZ'\sigma_u^2 = G\sigma_g^2$ , if  $\sigma_g^2 = N\sigma_u^2$

where  $N = \sum(p_i q_i)$

Therefore,  $G = ZDZ'/N$

$$D_{ii} = 1/(p_i q_i)$$

That is, assume the effect of SNPs is proportional to  $\sqrt{p_i q_i}$

So variance explained by SNPs is not affected by allele frequency

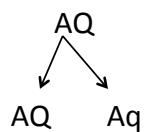
$$2c = 2a$$

17

## Estimate relationship from markers

Relationship depends on the markers or QTL

Eg QTL are due to recent mutations



Marker is the same but QTL is different

Rare SNP alleles tend to be a recent mutation

Therefore, treat SNPs differently according to MAF

18

## Estimate relationship from markers

Relationship depends on the markers or QTL  
Therefore, treat SNPs differently according to MAF

$$y = \text{mean} + g_1 + g_2 + g_3 + g_4 + g_5 + e$$

$$V(g_i) = (ZZ'/N)\sigma_i^2 \text{ for SNPs in MAF bin } i$$

19

## Estimate relationship from markers

Use haplotypes of markers

New definition of IBD for chromosome segments

Two segments are IBD if they coalesce without recombination  
A avoids definition of a base population

Chromosome segment homozygosity (CSH)

$$= P(\text{2 segments are IBD})$$

$$E(csh) = 1/(1+4N_e c)$$

20

## Estimate relationship from markers

$$E(csh) = 1/(1+4N_e c)$$

So CSH contains information about  $N_e$

If  $N_e$  varies in the past

Csh at small distances reflects  $N_e$  a long time ago

Csh at long distances reflects  $N_e$  a short time ago

CSH at  $c$  Morgans reflects  $N_e$  approximately  $1/(2c)$  generations ago

CSH at a range of distances contains information about past demography

21

## Estimate relationship from markers

CSH at a range of distances contains information about past demography

Problem: can't observe CSH directly

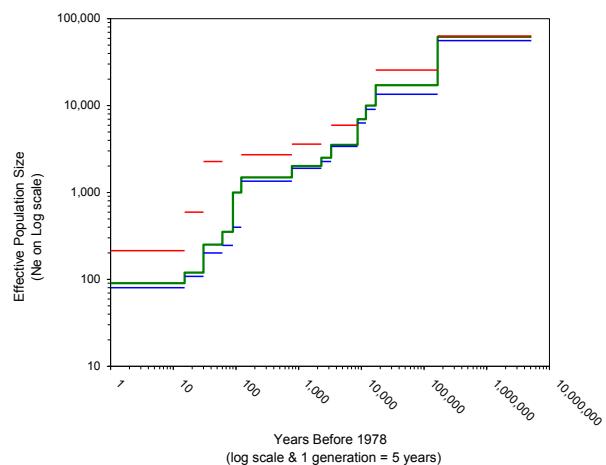
only observe haplotype homozygosity (HH)  
or runs of homozygosity (ROH)

Method to predict HH from past  $N_e$  (Macleod et al)

Therefore find past  $N_e$  that predicts observed HH

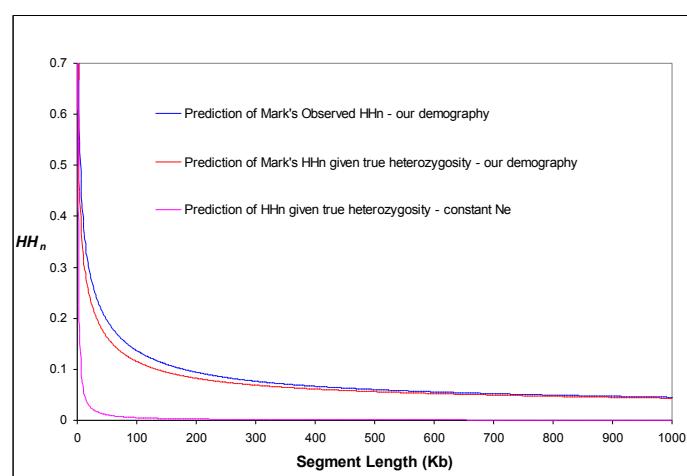
22

## Estimate relationship from markers



23

## Estimate relationship from markers



24

## Estimate relationship from markers

Can use HH or ROH in QTL mapping

### Additive effects

Calculate  $P(\text{QTL in position } x \text{ is IBD}) = P(\text{csh for surrounding chr})$

Eg  $P(\text{QTL IBD}) = 0.9$  if in middle of 10 identical markers

### Recessive effects

ROH within individual → homozygous QTL within the run

25

## Estimate relationship from markers

### Recessive effects

ROH within individual → homozygous QTL within the run

Detect embryonic lethals by missing ROH

26

# (Genome-wide) association analysis

Peter M. Visscher  
peter.visscher@uq.edu.au

1

## Key concepts

- Mapping QTL by association relies on linkage disequilibrium in the population;
- LD can be caused by close linkage between a QTL and marker (= good) or by confounding between a marker and other effects (= usually bad);
- Single SNP and GWAS analysis for quantitative traits follow the standard quantitative genetics model;
- Dense SNP data facilitate QC and testing for population structure;
- The power of QTL detection by LD depends on the proportion of phenotypic variance explained at a marker

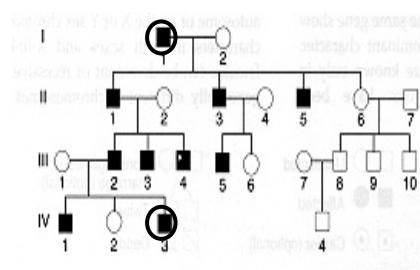
2

# Outline

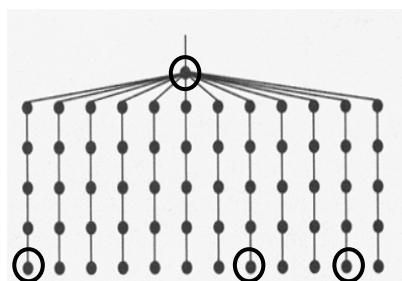
- Association vs linkage
- Linkage disequilibrium
- Analysis: single SNP
- GWAS: design, power
- GWAS: analysis

3

## Linkage



## Association

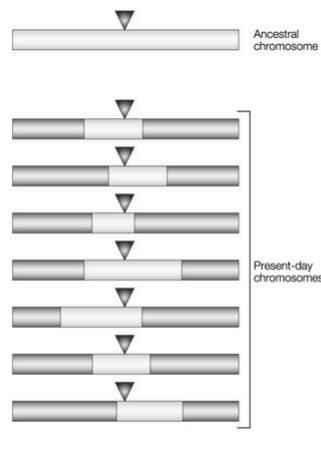


Families

Populations

4

## Linkage disequilibrium around an ancestral mutation



Nature Reviews | Genetics

5

[Ardlie et al. 2002]

## LD

- Non-random association between alleles at different loci
- Many possible causes
  - mutation
  - drift / inbreeding / founder effects
  - population stratification
  - selection
- Broken down by recombination

6

## Definition of D

- 2 bi-allelic loci
  - Locus 1, alleles A & a, with freq. p and (1-p)
  - Locus 2, alleles B & b with freq. q and (1-q)
  - Haplotype frequencies  $p_{AB}$ ,  $p_{Ab}$ ,  $p_{aB}$ ,  $p_{ab}$

$$D = p_{AB} - pq$$

7

$$r^2$$

$$r^2 = D^2 / [pq(1-p)(1-q)]$$

- Squared correlation between presence and absence of the alleles in the population
- ‘Nice’ statistical properties

8  
[Hill and Robertson 1968]

## Properties of $r$ and $r^2$

- Population in ‘equilibrium’  
 $E(r) = 0$   
 $E(r^2) = \text{var}(r) \approx 1/[1 + 4Nc] + 1/n$   
     $N$  = effective population size  
     $n$  = sample size (haplotypes)  
     $c$  = recombination rate
- $nr^2 \sim \chi_{(1)}^2$
- Human population is NOT in equilibrium

LD depends on  
population size and  
recombination  
distance

[Sved 1971; Weir and Hill 1980]<sup>9</sup>

## Population stratification

Both populations are in linkage equilibrium

	Allele frequency		Haplotype frequency			
	$p_{A1}$	$p_{B1}$	$p_{A1B1}$	$p_{A1B2}$	$p_{A2B1}$	$p_{A2B2}$
Pop. 1	0.9	0.9	0.81	0.09	0.09	0.01
Pop. 2	0.1	0.1	0.01	0.09	0.09	0.81
Average	0.5	0.5	0.41	0.09	0.09	0.41

Combined population:  $D = 0.16$  and  $r^2 = 0.41$

## Demonstrating stratification in a European American population

Catarina D Campbell<sup>1,2</sup>, Elizabeth L Ogburn<sup>1</sup>, Kathryn L Lunetta<sup>3,8</sup>, Helen N Lyon<sup>1,2</sup>, Matthew L Freedman<sup>4–6</sup>, Leif C Groop<sup>7</sup>, David Altshuler<sup>2,4,5</sup>, Kristin G Ardlie<sup>3</sup> & Joel N Hirschhorn<sup>1,2,4</sup>

**Table 2** No evidence for stratification using standard methods

SNPs	$\chi^2$ values <sup>a</sup>		Estimates of stratification parameters <sup>b</sup>			
	Median	Mean	$\lambda_{\max}$	$\lambda$	P	
Random SNPs	111	0.37	0.96	3.21	1	0.61
AIMs	67	0.58	0.95	—	—	0.61
Total	178	0.49	0.95	—	—	0.66

**Table 3** A strong association of *LCT*–13910C→T and height is reduced by rematching subjects on the basis of ancestry

N		Origin of grandparents <sup>c</sup>				
		All	Four US-born	Southeastern	Northwestern	
	Total	2,179	1,282	354	543	—
	Tall	1,123	645	127	351	—
	Short	1,056	637	227	192	—
<i>LCT</i> –13910 genotype counts <sup>e</sup>	Total	392,918,869	142,543,596	182,141,31	68,233,243	—
	Tall	16,147,4489	66,265,314	54,56,18	41,154,157	—
	Short	231,444,380	76,278,282	128,86,13	27,79,86	—
Hardy-Weinberg P	Total	$5.6 \times 10^{-7}$	0.5	0.89	0.89	—
	Tall	0.03	0.66	0.81	0.92	—
	Short	$2.5 \times 10^{-6}$	0.86	0.96	0.45	—
Association P		$3.6 \times 10^{-7}$	0.098	0.0016	0.71	0.0074
OR (95% c.i.) <sup>f</sup>		1.37 (1.22–1.54)	1.15 (0.97–1.36)	1.70 (1.22–2.38)	1.05 (0.81–1.37)	1.19 (1.05–1.36)

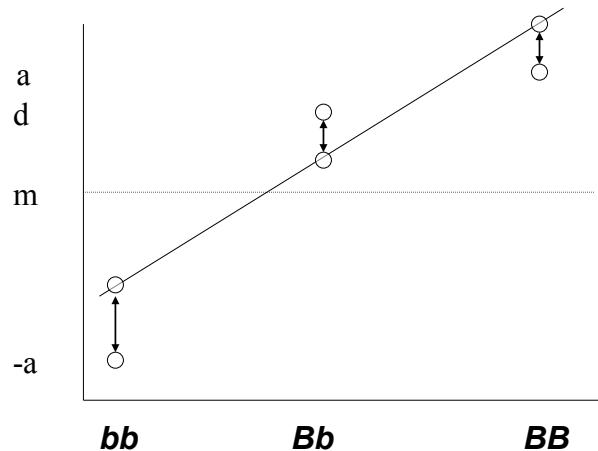
**Table 4** No association of *LCT*–13910C/T and height in other European populations

	Polish	Scandinavian	Combined
Genotypes (CC;CT;TT)	Tall 166,251,86 Short 174,235,96	— —	— —
Transmissions of T allele (TU) <sup>g</sup>	Tall — Short —	65,68 76,66	— —
P	0.92 0.99 (0.83–1.18)	0.43 0.91 (0.72–1.15)	0.58 0.96 (0.83–1.11)
OR (95% c.i.) <sup>h</sup>			11

## Analysis

- Single locus association
- TDT
- GWAS
- Least squares
- ML
- Bayesian methods

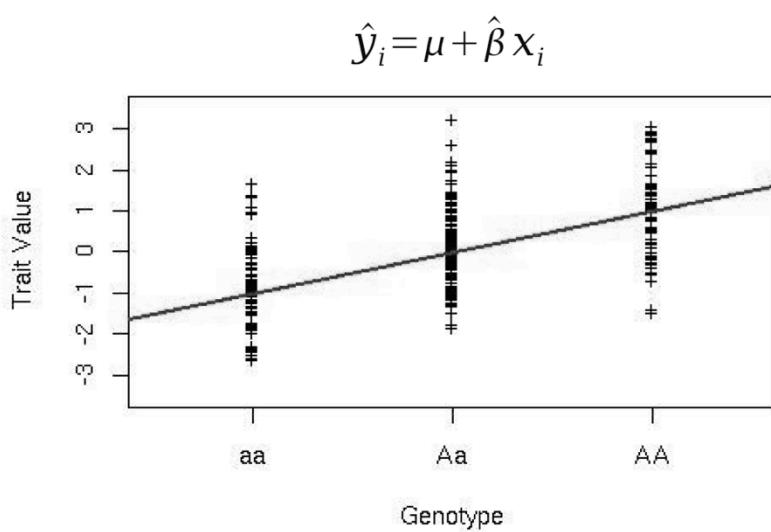
## Falconer model for single biallelic QTL



$$\begin{aligned}\text{Var}(X) &= \text{Regression Variance} + \text{Residual Variance} \\ &= \text{Additive Variance} + \text{Dominance Variance}\end{aligned}$$

13

## Unrelated Samples



14

## TDT

**= transmission disequilibrium test**

- Association with family-based controls
- Original TDT (disease mapping)
  - trios, two parents one affected progeny
  - test for transmission of allele to affected progeny from heterozygous parents
    - non-transmitted allele is the control
- **Quantitative traits**
  - Test for association between trait value and allele within parental mating type

15

## TDT for quantitative traits (regression model)

$$\hat{y}_{ij} = \mu + \hat{\beta}_b b_i + \hat{\beta}_w w_{ij}$$

$$b_{ij} = b_i = \begin{cases} \frac{g_{iF} + g_{iM}}{2} & \text{average of parental genotypes} \\ \sum_k \frac{g_{ik}}{n_{sibs}} & \text{average of sibling genotypes} \end{cases}$$

$$w_{ij} = g_{ij} - b_{ij}$$

$$g = 0, 1, 2$$

16

## Information for Within Test

- Families are only informative for the within family component when the offspring can have different genotypes....
  - AA x AA
  - AA x aa
  - **AA x Aa**
  - **Aa x Aa**
- At least one parent must be heterozygous

17

## Population Stratification

$$\hat{y}_{ij} = \mu + \hat{\beta}_b b_i + \hat{\beta}_w w_{ij}$$

- When there is no population stratification the slopes for the between and within test should be equal

18

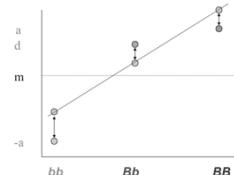
## Statistical power (linear regression)

$$y = \mu + \beta^*x + e, \quad x = 0, 1, 2$$

$$\sigma_p^2 = \sigma_q^2 + \sigma_e^2$$

$$\sigma_x^2 = 2p(1-p)$$

{HWE: note x is usually considered fixed in regression}



$$\sigma_q^2 = \beta^2 \sigma_x^2 = [a + d(1-2p)]^2 * 2p(1-p)$$

$$q^2 = \sigma_q^2 / \sigma_p^2$$

{QTL heritability}

19

## Statistical Power

$\chi^2$  test with 1 df:

$$E(X^2) = 1 + n R^2 / (1 - R^2)$$

$$= 1 + nq^2/(1-q^2)$$

$$= 1 + NCP$$

NCP = non-centrality parameter

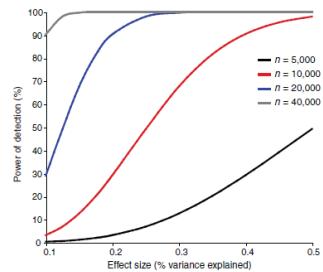
Power of association proportional to  $q^2$   
(Power of linkage proportional to  $q^4$ )

20

## Statistical Power ( $R$ )

```
alpha= 5e-8
threshold= qchisq(1-alpha,1)
q2= 0.005
n= 10000
ncp= n*q2/(1-q2)
power= 1-pchisq(threshold,1,ncp)
threshold
ncp
power
```

```
> alpha= 5e-8
> threshold= qchisq(1-alpha,1)
> q2= 0.005
> n= 10000
> ncp= n*q2/(1-q2)
> power= 1-pchisq(threshold,1,ncp)
> threshold
[1] 29.71679
> ncp
[1] 50.25126
> power
[1] 0.9492371
```



**Figure 1** Statistical power of detection in GWAS for variants that explain 0.1–0.5% of the variation at a type I error rate of  $5 \times 10^{-7}$  (calculated using the Genetic Power Calculator<sup>15</sup>). Shown is the power to detect a variant with a given effect size, assuming this type I error rate, which is typical for a GWAS with a sample size of  $n = 5,000\text{--}40,000$ .

21

## Power by association with SNP

(small effect; HWE)

$$\begin{aligned} \text{NCP(SNP)} &= n r^2 q^2 \\ &= r^2 * \text{NCP(causal variant)} \\ &= n * \{r^2 q^2\} = n * (\text{variance explained by SNP}) \end{aligned}$$

Power of LD mapping depends on the experimental sample size, variance explained by the causal variant and LD with a genotyped SNP

22

# Genetic Power Calculator (PGC)

<http://pngu.mgh.harvard.edu/~purcell/gpc/>

Genetic Power Calculator



## Genetic Power Calculator

S. Purcell & P. Sham, 2001-2009

This site provides automated power analysis for variance components (VC) quantitative trait locus (QTL) linkage and association tests in sibships, and other common tests. Suggestions, comments, etc to [Shaun Purcell](#).

If you use this site, please reference the following **Bioinformatics article**:

Purcell S, Cherny SS, Sham PC. (2003) Genetic Power Calculator: design of linkage and association genetic mapping studies of complex traits. *Bioinformatics*, 19(1):149-150.

Modules		Genetic Power Calculator	
<a href="#">Case-control for discrete traits</a>	<a href="#">Notes</a>	Total QTL variance :	<input type="text"/> (0 - 1)
<a href="#">Case-control for threshold-selected quantitative traits</a>	<a href="#">Notes</a>	Dominance : additive QTL effects :	<input type="text"/> (0 - 1) <input checked="" type="checkbox"/> No dominance (* see below)
<a href="#">QTL association for sibships and singletons</a>	<a href="#">Notes</a>	QTL increase allele frequency :	<input type="text"/> (0 - 1)
<a href="#">TDT for discrete traits</a>	<a href="#">Notes</a>	Marker M1 allele frequency :	<input type="text"/> (0 - 1)
<a href="#">TDT and parentTDT with ascertainment</a>	<a href="#">Notes</a>	Linkage disequilibrium (D-prime) :	<input type="text"/> (0 - 1)
<a href="#">TDT for threshold-selected quantitative traits</a>	<a href="#">Notes</a>	Sibling correlation :	<input type="text"/> (0 - 1) (* see below)
<a href="#">Epistasis power calculator</a>	<a href="#">Notes</a>	Sample Size :	<input type="text"/> (0 - 10000000) (N=families, not individuals)
<a href="#">QTL linkage for sibships</a>	<a href="#">Notes</a>	Sibship Size :	<input type="text"/> <input checked="" type="checkbox"/> Both parents genotyped
<a href="#">Probability Function Calculator</a>	<a href="#">Notes</a>	User-defined type I error rate :	<input type="text"/> 0.05 (0.0000001 - 0.5)
		User-defined power: determine N :	<input type="text"/> 0.80 (0 - 1) (1 - type II error rate)
		<input type="button" value="Process"/> <input type="button" value="Reset"/>	

# GWAS

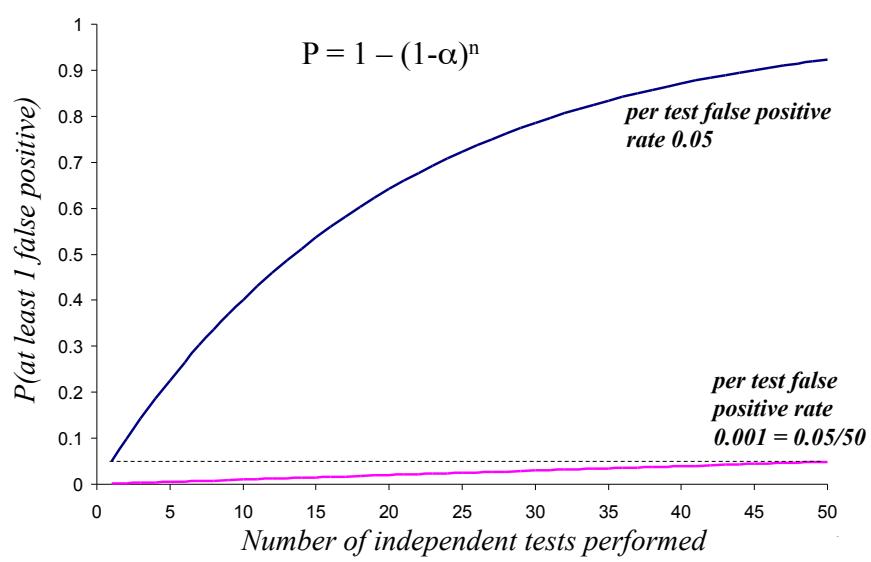
- Same principle as single locus association, but additional information
  - QC
    - Duplications, sample swaps, contamination
  - Power of multi-locus data
    - Unbiased genome-wide association
    - Relatedness
    - Population structure
    - Ancestry

## GWAS analysis

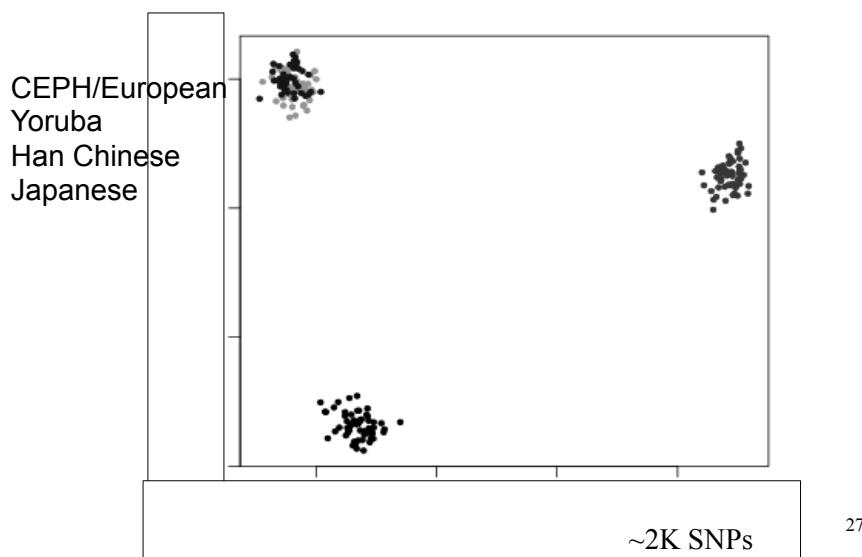
- Challenges
  - most obviously, multiple testing burden
  - computation
- Opportunities
  - simple methods can work well with ↑ data
  - novel analyses permitted, e.g.,
    - mixed model analyses
      - EMMAX, GCTA
    - enrichment/pathway analyses

25

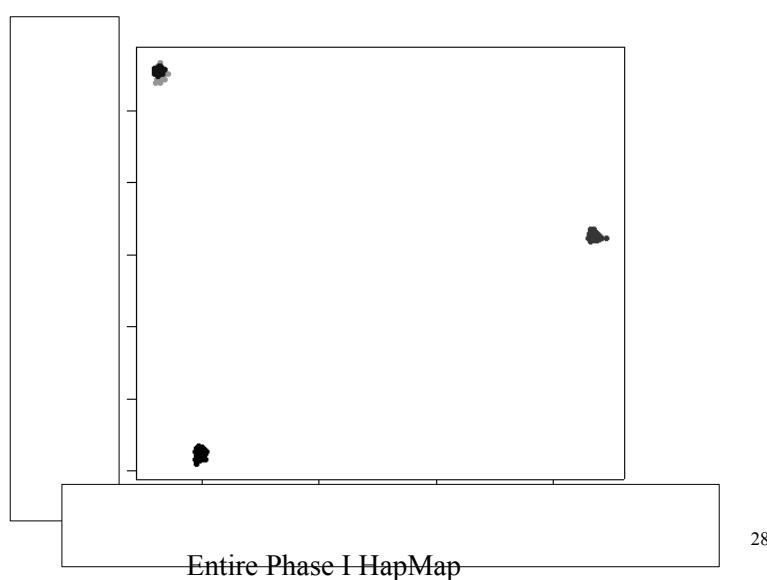
## The multiple testing burden



## Empirical assessment of ancestry

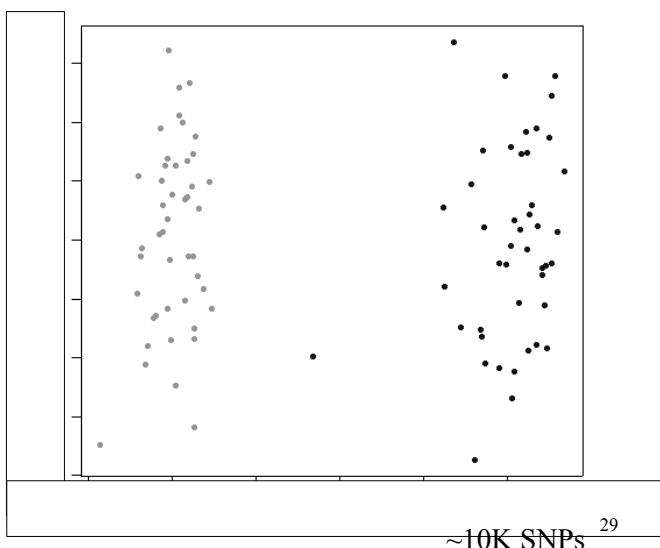


## Empirical assessment of ancestry



## Empirical assessment of ancestry

Han Chinese  
Japanese

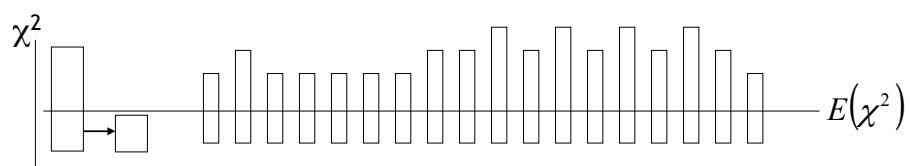
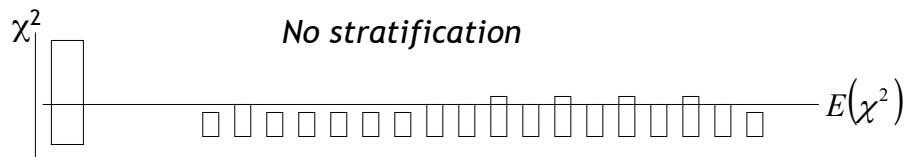


~10K SNPs <sup>29</sup>

## Many ways of dealing with structure

- Detect and discard ‘outliers’
- Detect, analysis and adjustment
  - E.g. genomic control
- Account for structure during analysis
  - Fit a few principal components as covariates
  - Fit entire SNP similarity matrix (= fitting all PCs)

## Genomic control



Stratification → adjust test statistic

31

## Genomic control

- Simple estimate of inflation factor

$$\hat{\lambda} = \text{median}\{\chi_1^2, \chi_2^2, \dots, \chi_N^2\}/0.456$$

```
> qchisq(0.5,1)
[1] 0.4549364
> qnorm(0.75)^2
[1] 0.4549364
```

- median protects from outliers
  - i.e. true effects
- bounded at minimum of 1
  - i.e. should never increase test statistic
- extends to multiple alleles, haplotypes, quantitative traits, different tests, etc

32

## Key concepts

- Mapping QTL by association relies on linkage disequilibrium in the population;
- LD can be caused by close linkage between a QTL and marker (= good) or by confounding between a marker and other effects (= usually bad);
- Single SNP and GWAS analysis for quantitative traits follow the standard quantitative genetics model;
- Dense SNP data facilitate QC and testing for population structure;
- The power of QTL detection by LD depends on the proportion of phenotypic variance explained at a marker

33

# Estimation of quantitative genetic parameters from distant relatives using marker data

Peter M. Visscher

peter.visscher@uq.edu.au

1

## Key concepts

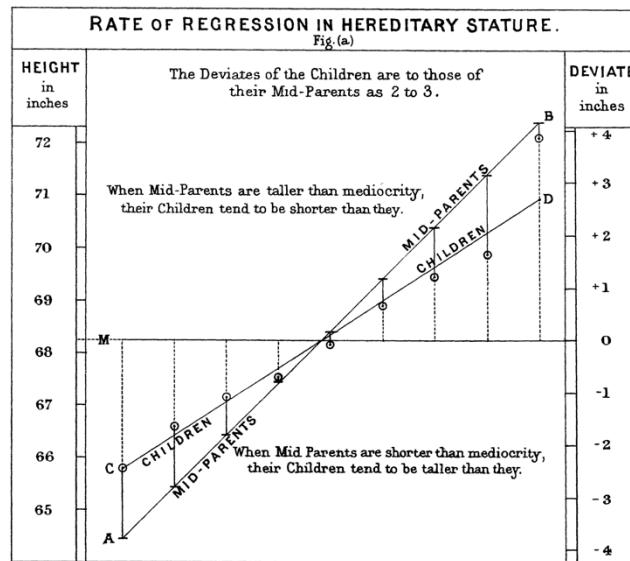
- Dense SNP panels allow the estimation of the expected genetic covariance between distant relatives ('unrelateds')
- A model based upon estimated relationships from SNPs is equivalent to a model fitting all SNPs simultaneously
- The total genetic variance due to LD between common SNPs and (unknown) causal variants can be estimated
- Genetic variance captured by common SNPs can be assigned to chromosomes and chromosome segments

2

ANTHROPOLOGICAL MISCELLANEA.

1886

REGRESSION towards MEDIOCRITY in HEREDITARY STATURE.  
By FRANCIS GALTON, F.R.S., &c.



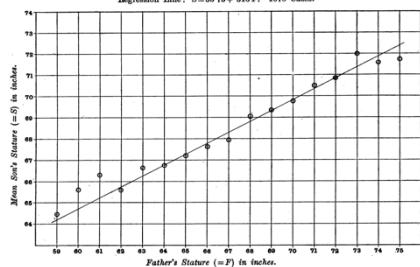
3

VOLUME II

NOVEMBER, 1903

No. 4

DIAGRAM I. Probable Stature of Son for given Father's Stature.  
Regression Line:  $S = 59.73 + .516 F$ . 1078 Cases.



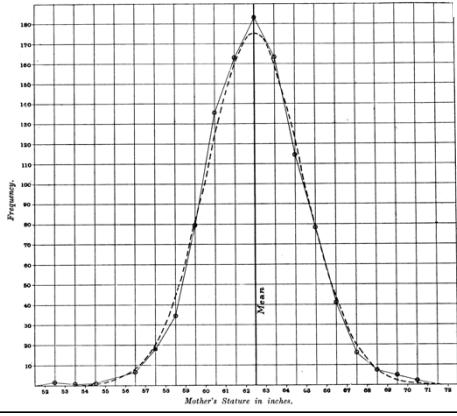
ON THE LAWS OF INHERITANCE IN MAN\*.

I. INHERITANCE OF PHYSICAL CHARACTERS.

By KARL PEARSON, F.R.S., assisted by ALICE LEE, D.Sc.  
University College, London.

364 *On the Laws of Inheritance in Man*

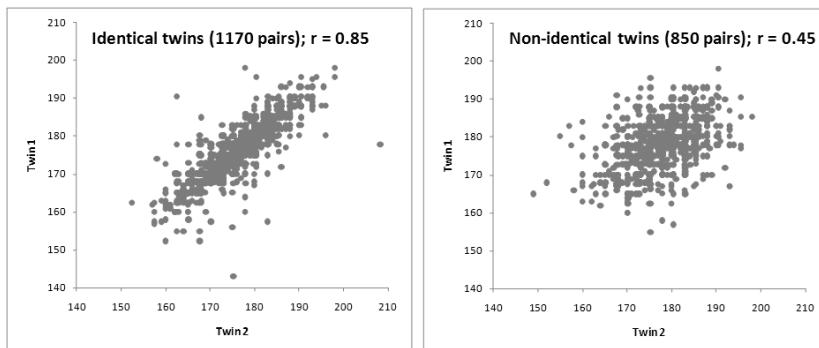
DIAGRAM IV. Distribution of Stature.



PAIR	CORRELATION	SE
Spouse	0.28	0.02
Son-Father	0.51	0.02
Daughter-Father	0.51	0.01
Son-Mother	0.49	0.02
Daughter-Mother	0.51	0.01
Brother-brother	0.51	0.03
Sister-sister	0.54	0.02
Brother-sister	0.55	0.01

4

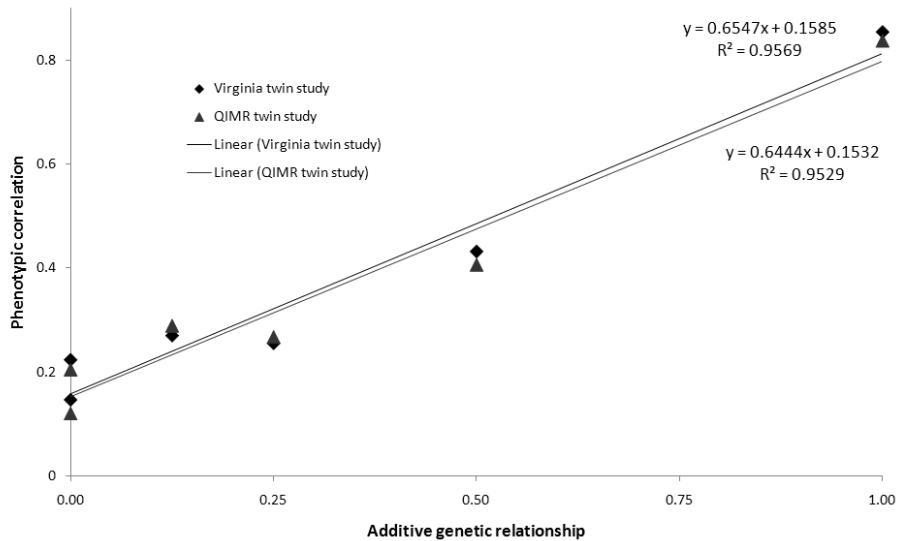
# 100 years later Heritability of human height



$$h^2 \sim 80\%$$

5

Based upon 1000s of twin families



6

PLOS ONE

Disease	Number of loci	Percent of Heritability Measure Explained	Heritability Measure
Age-related macular degeneration	5	50%	Sibling recurrence risk
Crohn's disease	32	20%	Genetic risk (liability)
Systemic lupus erythematosus	6	15%	Sibling recurrence risk
Type 2 diabetes	18	6%	Sibling recurrence risk
HDL cholesterol	7	5.2%	Phenotypic variance
Height	40	5%	Phenotypic variance
Early onset myocardial infarction	9	2.8%	Phenotypic variance
Fasting glucose	4	1.5%	Phenotypic variance

*open ACCESS freely available online*

**Rare Variants Create Synthetic Genome-Wide Associations**

Samuel P. Dickson<sup>1,2</sup>, Kai Wang<sup>3</sup>, Ian Krantz<sup>3,4,5</sup>, Hakon Hakonarson<sup>1,4,5</sup>, David B. Goldstein<sup>1\*</sup>

NATURE VOL 461 | November 2009

Where is the Dark Matter?

Vol 461 | 8 October 2009 | doi:10.1038/nature08494

nature

REVIEWS

Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hinds<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ransom<sup>8</sup>, Ian B. Cardon<sup>9</sup>, Aravinda Chakravarti<sup>10</sup>, Judy H. Cho<sup>11</sup>, Alan E. Guttmacher<sup>12</sup>, Augustine Kong<sup>13</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Meirik<sup>13</sup>, Charles N. Rotimi<sup>11</sup>, Montgomery Slatkin<sup>14</sup>, David Valle<sup>15</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

The case of the missing heritability

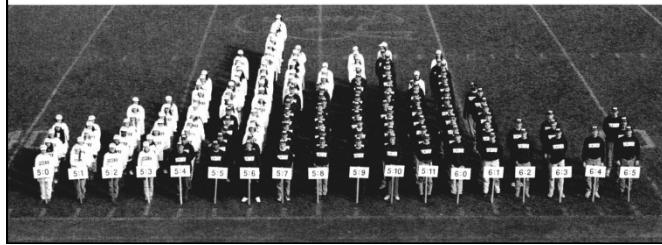
## Hypothesis testing vs. Estimation

- GWAS = hypothesis testing
  - Stringent p-value threshold
  - Estimates of effects biased (“Winner’s Curse”)
    - $E(bhat | test(bhat) > T) > b \{b \text{ fixed}\}$
    - $\text{var}(bhat) = \text{var}(b) + \text{var}(bhat | b) \{b \text{ random}\}$
- Can we estimate the total proportion of variation accounted for by all SNPs?

Common SNPs explain a large proportion of the heritability for human height

Jian Yang<sup>1</sup>, Beben Benyamin<sup>1</sup>, Brian P McEvoy<sup>1</sup>, Scott Gordon<sup>1</sup>, Anjali K Henders<sup>1</sup>, Dale R Nyholt<sup>1</sup>,  
Pamela A Madden<sup>2</sup>, Andrew C Heath<sup>2</sup>, Nicholas G Martin<sup>1</sup>, Grant W Montgomery<sup>1</sup>, Michael E Goddard<sup>3</sup> &  
Peter M Visscher<sup>1</sup>

## ***Are very distant relatives that share more of their genome by descent phenotypically more similar than those that share less?***



## Basic idea

- Estimates of additive genetic variance from known pedigree is unbiased
  - If model is correct
  - Despite variation in identity given the pedigree
  - Pedigree gives correct expected IBD
- Unknown pedigree: estimate genome-wide IBD from marker data
  - Estimate additive genetic variance given this estimate of relatedness
- Idea is not new
  - (Evolutionary) genetics literature (Ritland, Lynch, Hill, others)

## Close vs distant relatives

- Detection of close relatives (fullsibs, parent-offspring, halfsibs) from marker data is relatively straightforward
- But close relatives may share environmental factors
  - Biased estimates of genetic variance
- Solution: use only (very) distant relatives

11

## A model for a single causal variant

	AA	AB	BB
frequency	$(1-p)^2$	$2p(1-p)$	$p^2$
x	0	1	2
effect	0	b	2b
$z = [x - E(x)]/\sigma_x$	$-2p/\sqrt{2p(1-p)}$	$(1-p)/\sqrt{2p(1-p)}$	$2(1-p)/\sqrt{2p(1-p)}$

$$y_j = \mu' + x_{ij}b_i + e_j \quad x = 0, 1, 2 \text{ {standard association model}}$$

$$y_j = \mu + z_{ij}u_j + e_j \quad u = b\sigma_x; \mu = \mu' + b\sigma_x$$

12

## Multiple (m) causal variants

$$y_j = \mu + \sum z_{ij} u_j + e_j$$

$$= \mu + g_j + e_j$$

$$\mathbf{y} = \mu \mathbf{1} + \mathbf{g} + \mathbf{e}$$

$$= \mu \mathbf{1} + \mathbf{Zu} + \mathbf{e}$$

13

## Equivalence

Let  $u$  be a random variable,  $u \sim N(0, \sigma_u^2)$

Then  $\sigma_g^2 = m\sigma_u^2$  and

$$\begin{aligned}\text{var}(\mathbf{y}) &= \mathbf{Z}\mathbf{Z}' \sigma_u^2 + \mathbf{I}\sigma_e^2 \\ &= \mathbf{Z}\mathbf{Z}' (\sigma_g^2/m) + \mathbf{I}\sigma_e^2 \\ &= \mathbf{G} \sigma_g^2 + \mathbf{I}\sigma_e^2\end{aligned}$$

Model with individual genome-wide additive values using relationships ( $\mathbf{G}$ ) at the causal variants is equivalent to a model fitting all causal variants

We can estimate genetic variance just as if we would do using pedigree relationships

## But we don't have the causal variants

If we estimate  $\mathbf{G}$  from SNPs:

- lose information due to imperfect LD between SNPs and causal variants
- how much we lose depends on
  - density of SNPs
  - allele frequency spectrum of SNPs vs. causal variants
- estimate of variance → missing heritability

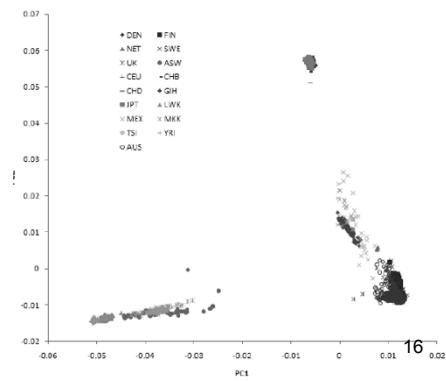
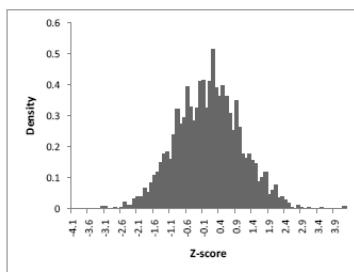
Let  $\mathbf{A}$  be the estimate of  $\mathbf{G}$  from N SNPs:

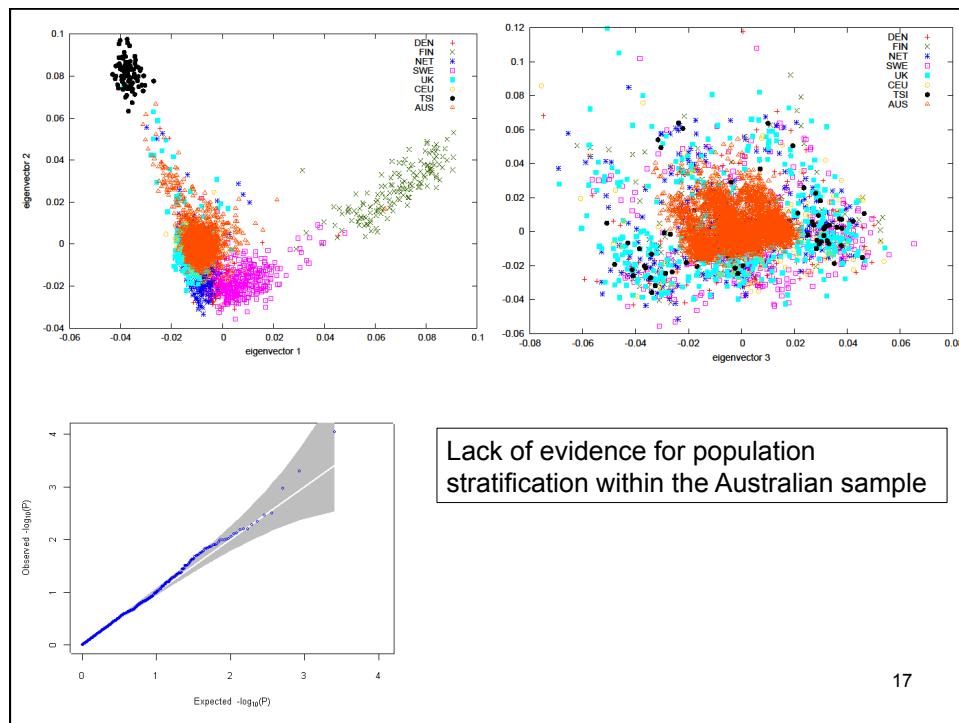
$$\begin{aligned} A_{jk} &= (1/N) \sum \{ x_{ij} - 2p_i \} (x_{ik} - 2p_i) / \{ 2p_i(1-p_i) \} \\ &= (1/N) \sum z_{ij} z_{ik} \end{aligned}$$

15

## Data

- ~4000 ‘unrelated’ individuals
- Ancestry ~British Isles
- Measurement on height (self-report or clinically measured)
- GWAS on 300k (‘adults’) or 600k (16-year olds) SNPs





## Methods

- Estimate realised relationship matrix from SNPs  $y_i = g_i + e_i$   $\text{var}(y) = V = A\sigma_g^2 + I\sigma_e^2$
- Estimate additive genetic variance

$$A_{ijk} = \frac{\text{cov}(x_{ij}a_i, x_{ik}a_i)}{\sqrt{\text{var}(x_{ij}a_i)\text{var}(x_{ik}a_i)}} = \frac{\text{cov}(x_{ij}, x_{ik})}{2p_i(1-p_i)}$$

Base population = current population

$$A_{jk} = \frac{1}{N} \sum_i A_{ijk} = \begin{cases} \frac{1}{N} \sum_i \frac{(x_{ij} - 2p_i)(x_{ik} - 2p_i)}{2p_i(1-p_i)}, & j \neq k \\ 1 + \frac{1}{N} \sum_i \frac{x_{ij}^2 - (1+2p_i)x_{ij} + 2p_i^2}{2p_i(1-p_i)}, & j = k \end{cases}$$

18

# Statistical analysis

$$\text{var}(\mathbf{y}) = \mathbf{V} = \mathbf{A}\sigma_g^2 + \mathbf{I}\sigma_e^2$$

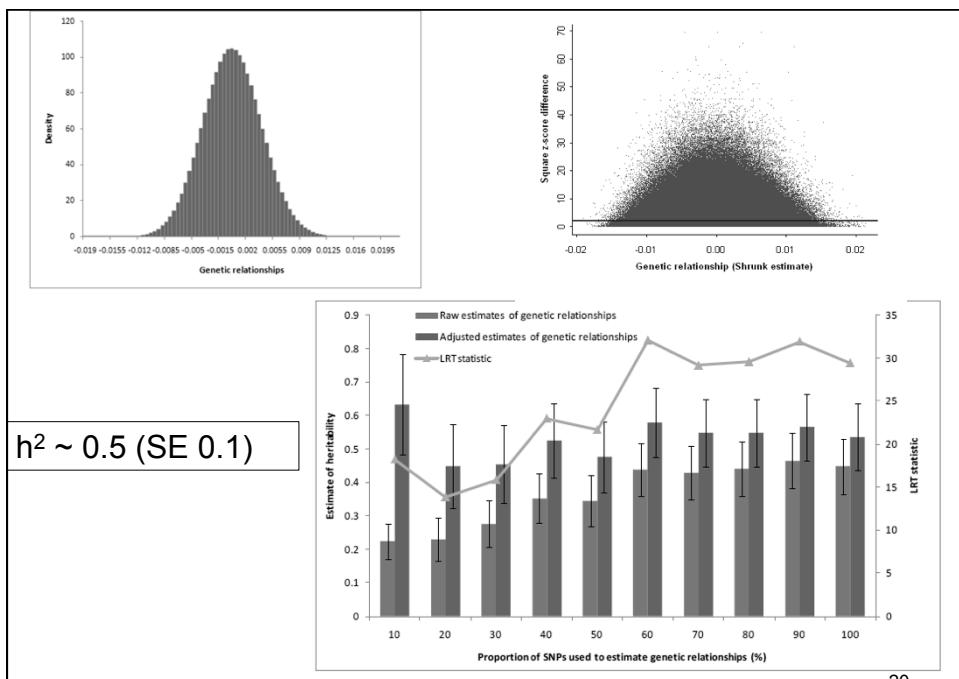
$\mathbf{y}$  standardised  $\sim N(0,1)$

No fixed effects other than mean

$\mathbf{A}$  estimated from SNPs

Residual maximum likelihood (REML)

19



## Checking for population structure

**Table 1**  
Estimates of the Variance Explained by the SNPs on Even Chromosomes from 10 Simulation Replicates

Replicate	$h^2$	SE
1	0.045	0.055
2	0.025	0.057
3	0.0	0.058
4	0.0	0.057
5	0.0	0.059
6	0.0	0.056
7	0.057	0.056
8	0.0	0.062
9	0.0	0.057
10	0.0	0.054

Note: A total of 1,000 causal variants were simulated on the odd chromosomes, with a total heritability of 0.8. Genetic variance was estimated from a relationship matrix constructed from all SNPs on the even chromosomes. The same genotypes were used as in Yang et al. (2010). If there is population structure then estimated relatedness on the even chromosomes is correlated with relatedness on the odd chromosomes (where the causal variants are simulated) and therefore genetic variance will be associated with the even chromosomes.

21

## Partitioning variation

- If we can estimate the variance captured by SNPs genome-wide, we should be able to partition it and attribute variance to regions of the genome
- “Population based linkage analysis”

22

# Genome partitioning

- Partition additive genetic variance according to groups of SNPs
  - Chromosomes
  - Chromosome segments
  - MAF bins
  - Genic vs non-genic regions
  - Etc.
- Estimate genetic relationship matrix from SNP groups
- Analyse phenotypes by fitting multiple relationship matrices
- Linear model & REML (restricted maximum likelihood)

## REPORT

### GCTA: A Tool for Genome-wide Complex Trait Analysis

Jian Yang,<sup>1,\*</sup> S. Hong Lee,<sup>1</sup> Michael E. Goddard,<sup>2,3</sup> and Peter M. Visscher<sup>1</sup>

# Data from the GENEVA Consortium

- Investigators: Bruce Weir, Teri Manolio and many others
- Data
  - ~14,000 European Americans
    - ARIC
    - NHS
    - HPFS
  - Affy 6.0 genotype data
    - ~600,000 after stringent QC
  - Phenotypes on height, BMI, vWF and QT Interval

### Genome partitioning of genetic variation for complex traits using common SNPs

Jian Yang<sup>1\*</sup>, Teri A Manolio<sup>2</sup>, Louis R Pasquale<sup>3</sup>, Eric Boerwinkle<sup>4</sup>, Neil Caporaso<sup>5</sup>, Julie M Cunningham<sup>6</sup>, Mariza de Andrade<sup>7</sup>, Bjarke Feenstra<sup>8</sup>, Eleanor Feingold<sup>9</sup>, M Geoffrey Hayes<sup>10</sup>, William G Hill<sup>11</sup>, Maria Teresa Landi<sup>12</sup>, Alvaro Alonso<sup>13</sup>, Guillaume Lettre<sup>14</sup>, Peng Lin<sup>15</sup>, Hua Ling<sup>16</sup>, William Lowe<sup>17</sup>, Rasika A Mathias<sup>18</sup>, Mads Melbye<sup>8</sup>, Elizabeth Pugh<sup>16</sup>, Marilyn C Cornelis<sup>19</sup>, Bruce S Weir<sup>20</sup>, Michael E Goddard<sup>21,22</sup> & Peter M Visscher<sup>1</sup>

## QC of SNPs

Table 9. Summary of recommended SNP filters. "Broad" refers to SNPs failed by the genotyping center and "CC" refers to filters recommended by the GENEVA Coordinating Center.

SNP's kept	SNP's lost	remove SNPs with:
909,622	0	
843,985	65,637	Broad: call rate < 95%
841,820	2,165	Broad: plate associations (>6 plates with p<1e-10)
		CC: one member of each pair of duplicate probes (mostly AFFX probes)
839,046	2,774	CC: MAF = 0 in all samples
838,715	331	CC: call rate < 95%
838,493	222	CC: >5 discordant calls in 307 pairs of duplicates
802,026	36,468	CC: sex difference in allelic frequency between sexes > 0.10 in either European- or African-ancestry groups
801,956	0	CC: sex difference in heterozygosity > 0.3 in either ancestry group (for autosomal or XY)
801,956	0	CC: Hardy-Weinberg p-value < 1e-3 in either European- or African ancestry group
780,062	21,894	

- 780,062 SNPs after QC steps listed in the table.
- Exclude 141,772 SNPs with MAF < 0.02 in European-ancestry group.
- Exclude 36,949 SNPs with missingness > 2% in all samples.
- Include autosomal SNPs only.
- End up with 577,778 SNPs.

25

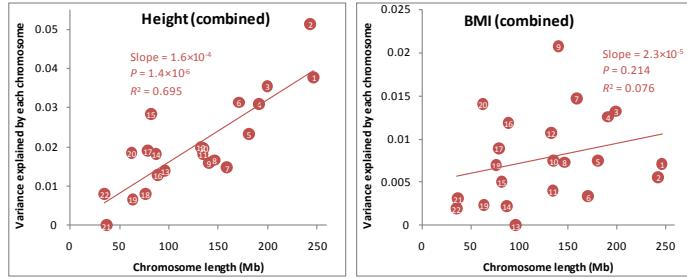
## Results (genome-wide)

Table 1 Estimates of the variance explained by all autosomal SNPs for height, BMI, vWF and QT<sub>i</sub>

Trait	<i>n</i>	No PCs <sup>a</sup>		10 PCs <sup>b</sup>		Heritability <sup>d</sup>	GWAS <sup>e</sup>
		$h_G^2$ (s.e.) <sup>c</sup>	<i>P</i>	$h_G^2$ (s.e.)	<i>P</i>		
Height	11,576	0.448 (0.029)	$4.5 \times 10^{-69}$	0.419 (0.030)	$7.9 \times 10^{-48}$	80–90% <sup>32</sup>	~10% <sup>23</sup>
BMI	11,558	0.165 (0.029)	$3.0 \times 10^{-10}$	0.159 (0.029)	$5.3 \times 10^{-9}$	42–80% <sup>25,26</sup>	~1.5% <sup>14</sup>
vWF	6,641	0.252 (0.051)	$1.6 \times 10^{-7}$	0.254 (0.051)	$2.0 \times 10^{-7}$	66–75% <sup>33,34</sup>	~13% <sup>15</sup>
QT <sub>i</sub>	6,567	0.209 (0.050)	$3.1 \times 10^{-6}$	0.168 (0.052)	$5.0 \times 10^{-4}$	37–60% <sup>35,36</sup>	~7% <sup>16</sup>

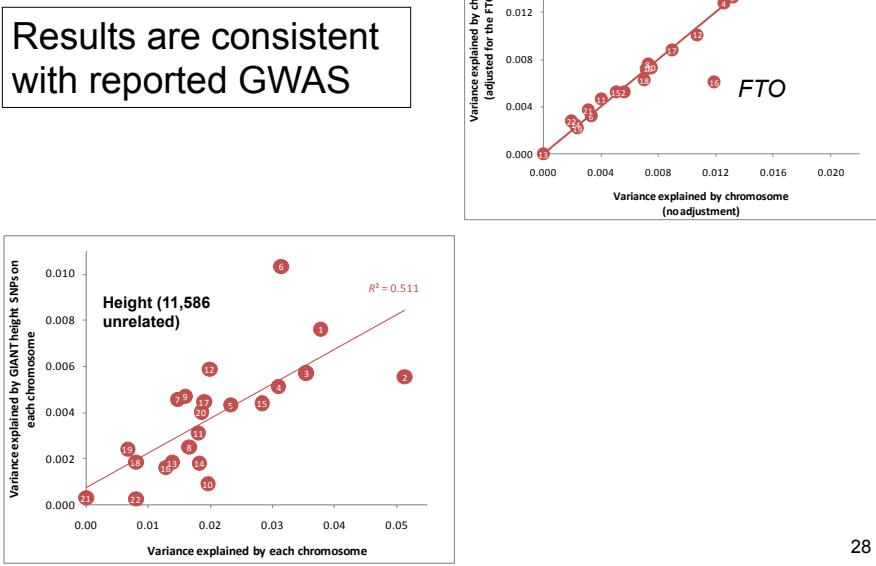
26

## Genome-partitioning: longer chromosomes explain more variation



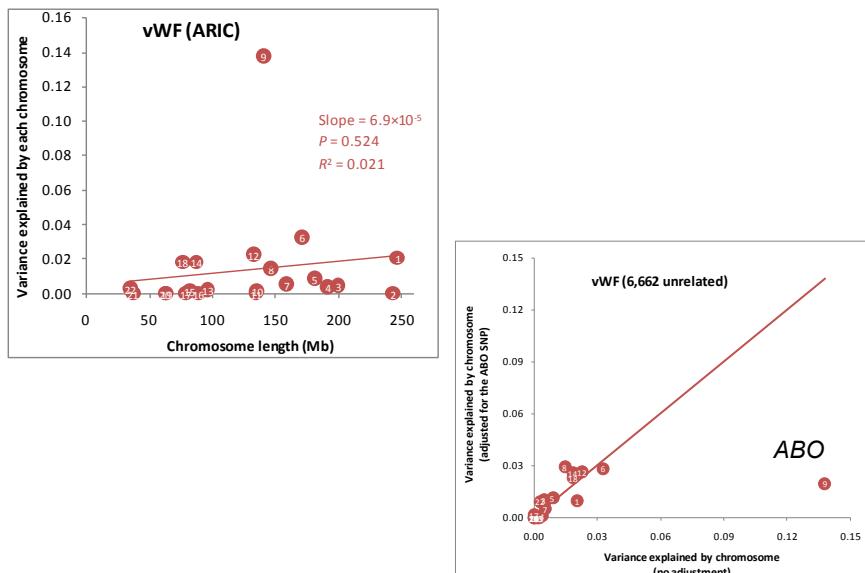
27

## Results are consistent with reported GWAS

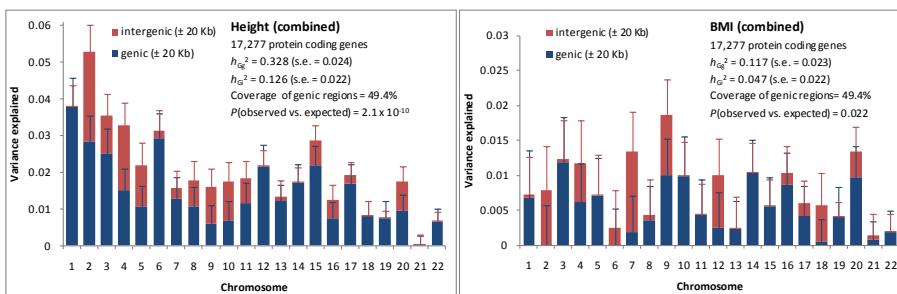


28

## Inference robust with respect to genetic architecture



## Genic regions explain variation disproportionately



## Key concepts

- Dense SNP panels allow the estimation of the expected genetic covariance between distant relatives ('unrelateds')
- A model based upon estimated relationships from SNPs is equivalent to a model fitting all SNPs simultaneously
- The total genetic variance due to LD between common SNPs and (unknown) causal variants can be estimated
- Genetic variance captured by common SNPs can be assigned to chromosomes and chromosome segments

# Prediction of quantitative traits using marker data

Peter M. Visscher & Michael E. Goddard

peter.visscher@uq.edu.au

Mike.Goddard@depi.vic.gov.au

1

## Key concepts

- Prediction of phenotypic values is limited by heritability
- Accuracy of prediction depends on
  - how well marker effects are estimated (sample size)
  - how well marker effects are correlated with causal variants (LD)
- Estimation of marker effects and prediction in the same data leads to (severe) bias
  - winner's curse; over-fitting
- Variance explained by a SNP-based predictor is not the same as the variance explained by those SNPs
- Marker data captures both between and within family genetic variation
- Best prediction methods take genetic values as random effects

2

 **WORLD WIDE SIRES, LTD.**  
 YOUR FOUNDATION...YOUR FUTURE

**7H010780 UNICORN MILLION ABERLIN-ET \*TR \*TV**  
 \*TL \*TY \*TD  
 USA 000066985571  
**MILLION X GOLDYN X O MAN**  
 100% Registered Holstein Ancestry



**ABERLIN**



*Copyright © 2001 by the Genetics Society of America*  
*Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps*  
*T. H. E. Meuwissen,<sup>a</sup> B. J. Hayes<sup>b</sup> and M. E. Goddard<sup>a,c</sup>*

Production	PTA	Management Traits
TPI	2206	3.32 SCE / Rel1.%
NMS	561	3.53 DCE / Rel1.%
PTA Milk (lbs)	836	2.53 SSB / Rel1.%
PTA Protein (lbs)	29	2.31 DSB / Rel1.%
PTA Protein (%)	0.01	1.90 SCS
PTA Fat (lbs)	38	72 Productive Life
PTA Fat (%)	0.02	0.0 DPR / Rel1.%
Production Reliability %	76	24 CCP / Rel1.%
Dtrs / Herds	aAa	
	0/0	

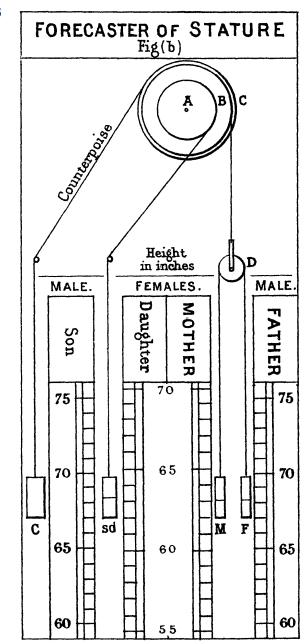
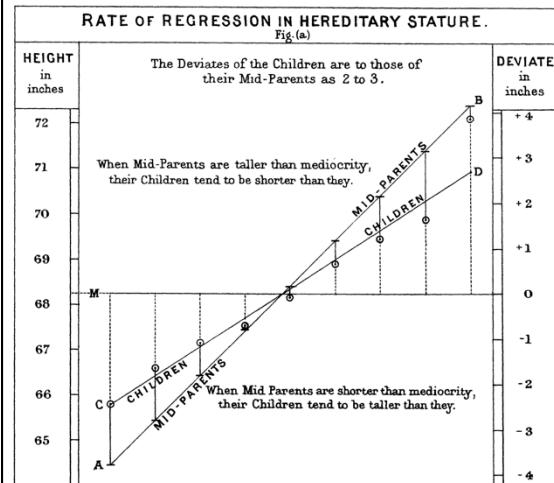
**DAM RC-LC Goldyn ATM**

"Genomic selection" = individual prediction in a commercial setting

## Take-home from animal breeding

- (1) Don't need genome-wide significant effects
- (2) Don't need to know causal variants
- (3) Don't need to know function

Regression Towards Mediocrity in Hereditary Stature.  
 Author(s): Francis Galton  
 Source: *The Journal of the Anthropological Institute of Great Britain and Ireland*, Vol. 15 (1886), pp. 246-263



## A quantitative genetics model

$$y = \text{fixed effects} + G + E$$

$$G = A + D + I$$

Possible predictions:

- Predict  $y$  from fixed effects and  $G$
- Predict  $G$  from  $A$
- Predict  $y$  from  $A$
- **Predict  $y$  from  $A$  using markers**

## Prediction using linear regression

$$y = \beta * x + e$$

- Usually,  $b$  and  $x$  are considered ‘fixed’
- For SNPs,  $x$  is random with variance  $2p(1-p)$  assuming HWE
- Later we will consider the case where  $\beta$  is random

7

## Theory (additive model) $m$ unlinked causal variants

$$y_i = \sum_{j=1}^m x_{ij} b_j + e_i = a_i + e_i$$

$$\text{var}(y) = \sum_{j=1}^m \text{var}(x_j) b_j^2 + \text{var}(e) = \text{var}(a) + \text{var}(e)$$

$$\begin{aligned}\text{cov}(y_i, y_k) &= \sum_{j=1}^m \text{cov}(x_{ij}, x_{kj}) b_j^2 + \text{cov}(e_i, e_k) \\ &= \text{cov}(a_i, a_k) + \text{cov}(e_i, e_k) \\ &= \text{cov}(a_i, a_k) \text{ if } \text{cov}(e_i, e_k) = 0\end{aligned}$$

8

## Prediction

$$\hat{y}_i = \sum_{j=1}^m x_{ij} \hat{b}_j = \hat{a}_i$$

$$\text{var}(\hat{y}) = \sum_{j=1}^m \text{var}(x_j) \hat{b}_j^2 = \text{var}(\hat{a})$$

$$\text{cov}(\hat{y}_i, \hat{y}_k) = \sum_{j=1}^m \text{cov}(x_{ij}, x_{kj}) \hat{b}_j^2 = \text{cov}(\hat{a}_i, \hat{a}_k)$$

9

## - theory -

$$\begin{aligned} \text{cov}(\hat{y}_i, y_i) &= \text{cov}\left\{\sum_{j=1}^m (x_{ij} \hat{b}_j), \sum_{j=1}^m x_{ij} b_j + e_i\right\} \\ &= \sum_{j=1}^m \text{var}(x_{ij}) \hat{b}_j b_j + \sum_{j=1}^m x_{ij} \text{cov}(\hat{b}_j, e_i) \end{aligned}$$

If  $b$  estimated from the same data in which prediction is made, then the second term is non-zero

10

## - theory -

$m$  markers, sample size  $N$

All  $b = 0$

Multiple linear regression of  $y$  on  $m$  markers

$$E(R^2) = m/N \quad \{ \text{strictly } m/(N-1) \}$$

→ Variation “explained” by chance

[Wishart, 1931]

11

## Selection bias

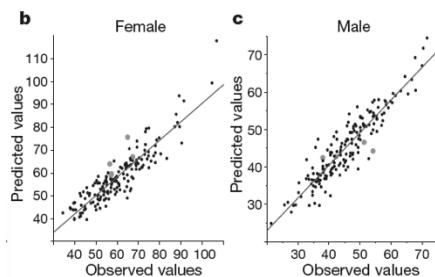
- Select  $m$  ‘best’ markers out of  $M$  in total
- ‘Prediction’ in same sample

$$E(R^2) >> m/N$$

→ Lots of variation explained by chance

ARTICLE

The *Drosophila melanogaster*  
Genetic Reference Panel



~15 best markers selected from 2.5 million markers

12

## Effect of errors in estimating SNP effects (least squares; single SNP)

$$y_i = x_i b + e_i$$

$$\hat{b} = b + \varepsilon$$

$$E(\hat{b}) = b$$

$$\text{var}(\hat{b}) = \text{var}(\varepsilon) = \sigma_e^2 / \sum x^2 \approx \text{var}(y) / \{N \text{var}(x)\}$$

$\text{var}(x) = 2p(1-p)$  under HWE

Define  $R_{SNP}^2 = \text{var}(x)b^2 / \text{var}(y)$

= contribution of single SNP to heritability

13

## - effects of errors -

$$\hat{R}_{y,\hat{y}}^2 = \text{cov}(y, \hat{y})^2 / \{\text{var}(y) \text{var}(\hat{y})\}$$

$$\begin{aligned} E[\text{cov}(y, \hat{y})] &= E[\text{cov}(xb, x\hat{b})] = \text{var}(x_i)E(\hat{b})b \\ &= \text{var}(x)b^2 \end{aligned}$$

$$\begin{aligned} E[\text{var}(\hat{y})] &= E[\text{var}(x\hat{b})] = \text{var}(x)E[\hat{b}^2] \\ &= \text{var}(x)[b^2 + \text{var}(\hat{b})] \approx \text{var}(x)b^2 + \text{var}(x)\text{var}(y)/[N \text{var}(x)] \\ &= \text{var}(x)b^2 + \text{var}(y)/N \end{aligned}$$

$$E(\hat{R}_{y,\hat{y}}^2) \approx R_{SNP}^2 / [1 + 1/\{NR_{SNP}^2\}]$$

14

## *m* variants

$$R_m^2 = \text{var}(a) / \text{var}(y) = h^2$$

$$E(\hat{R}_{y,\hat{y}}^2) \approx h^2 / [1 + m / \{Nh^2\}]$$

Even if we knew all  $m$  causal variants but needed to estimate their effect sizes then the variance explained by the predictor is less than the variance explained by the causal variants in the population.

[Daetwyler et al. 2008, PLoS Genetics; Visscher, Yang, Goddard 2010, Twin Research Human Genetics 2010]<sup>15</sup>

## Take-home

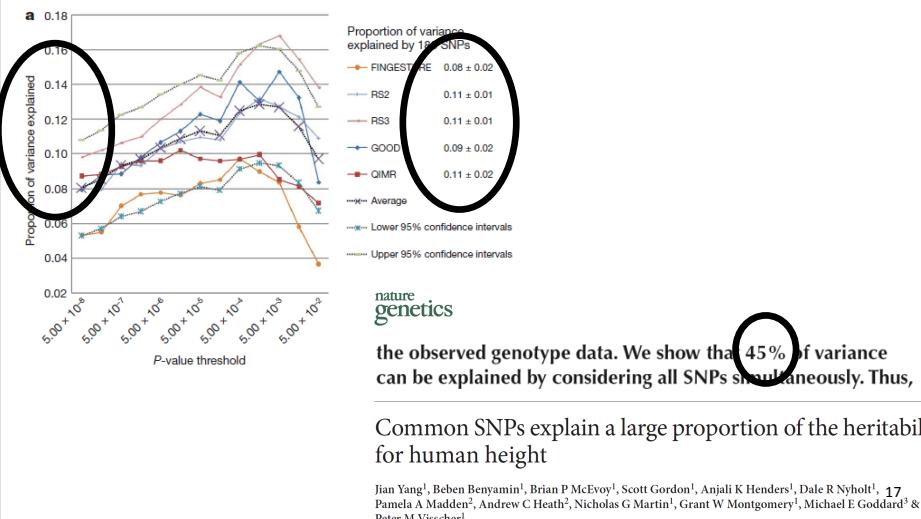
(4) Estimation of variance contributed by (all) loci is not the same as prediction accuracy

unless the effect sizes are estimated without error

## LETTER

doi:10.1038/nature09410

### Hundreds of variants clustered in genomic loci and biological pathways affect human height



### Measures of how well a predictor works

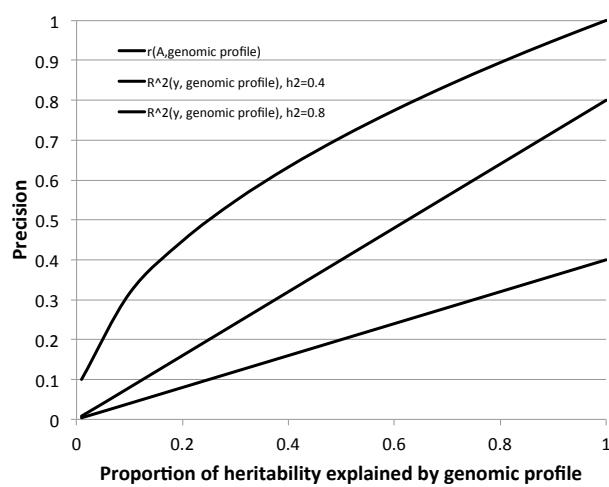
- “Accuracy” (animal breeding)
  - Correlation between true genome-wide genetic value and its predictor
- $R^2$  from a regression of outcome on predictor (human genetics)
- Area-under-curve from ROC analyses (disease classification)

## Limits of prediction

- A perfect predictor of A can be a lousy predictor of a phenotype
- The regression  $R^2$  has a maximum that depends on heritability
- The regression  $R^2$  is limited by unknown (eg future) fixed effects and covariates

19

## Predictions from known variants



20

### Prediction using genetic markers: using between and within-family genetic variation

FAMILY HISTORY INDIVIDUAL GENETIC RISK



All members of a sibship have equal predicted risk  
Between family variance

Members of a sibship have individual predicted risk  
Between and within family variance

21

### In class demo

- 180 height variants from Lango-Allan et al. 2010
  - Estimation of  $b$  from data ( $N \sim 4000$ )
    - Note that  $E(R^2) = 180/4000 = 0.045$  by chance!
  - Using  $b$  from Lango-Allen paper
- Taking the top 180 SNPs from GWAS

22

## Analysis demonstration

- **Data:**
  - Genotype data: 3,924 unrelated individuals and ~2.5M SNPs.
  - Phenotype data: height z-scores (adjusted for age and sex)
  - 180 SNPs identified by the GIANT meta-analysis (MA) of height ( $n = \sim 180,000$ )
- **Analyses:**
  - Estimating effect sizes of the 180 height SNPs in the data.
  - PLINK scoring: 180 GIANT SNPs, using effect sizes estimated from GIANT MA.
  - GWAS analysis in the data, selecting top SNPs at 180 loci and predicting the phenotypes in the same data.
- **Results:**
  - Estimation:  $R^2 = 0.134$  ( $R^2 = 0.046$  by chance), adjusted  $R^2 = 0.093$
  - Prediction:  $R^2 = 0.09$
  - Prediction using the top SNPs selected in the same data:  $R^2 = 0.429$

23

## Identifying people at high risk: T1D

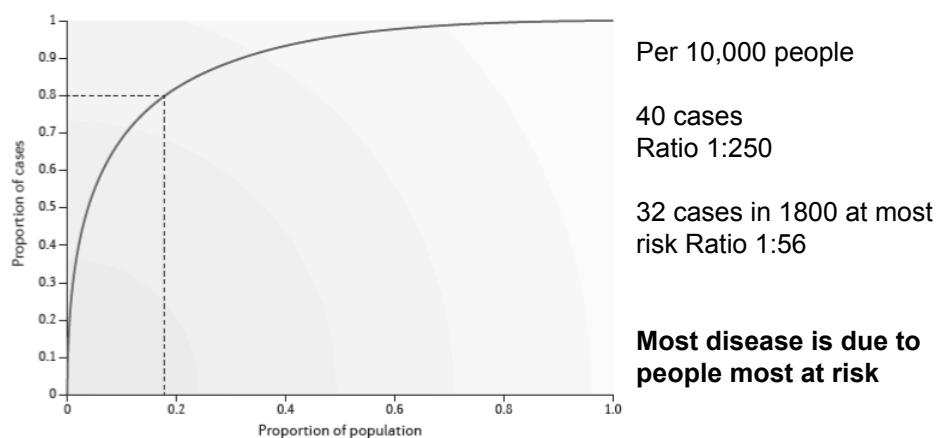


Figure 3 | The receiver operating characteristic (ROC) curve for the known T1D loci. The ROC curve plots the sensitivity of genetic type 1 diabetes (T1D) prediction

Polychronakos & Li NRG 2011  
Clayton PLoS Genetics 2009

## Prediction of genetic value using better predictors

Model with additive inheritance

$$y = g + e$$

$$V(g) = G\sigma_g^2, V(e) = I\sigma_e^2, V(y) = V = G\sigma_g^2 + I\sigma_e^2,$$

Aim is to predict  $g$  for individuals

Eg to predict future risk of a disease

25

## Prediction of genetic value

$$y = g + e$$

$$V(g) = G\sigma_g^2, V(e) = I\sigma_e^2, V(y) = V = G\sigma_g^2 + I\sigma_e^2,$$

Best prediction is

$$\hat{g} = E(g | y)$$

If  $y$  and  $g$  are bivariate normal

$$E(g | y) = b'y = \sigma_g^2 G V^{-1} y$$

26

## Prediction of genetic value

Eg Unrelated individuals

$$V(g) = I h^2, V(e) = I(1-h^2), V(y) = I,$$

Best prediction is

$$\hat{g} = E(g | y) = b'y = \sigma_g^2 G V^{-1} y = h^2 y$$

27

## Prediction of genetic value

$$y = g + e, g = Zu$$

$$V(u) = I \sigma_u^2, V(Zu) = ZZ' \sigma_u^2,$$

Best prediction is

$$\hat{u} = E(u | y)$$

If  $y$  and  $u$  are multivariate normal

$$E(u | y) = b'y = \sigma_u^2 Z' V^{-1} y$$

28

## Prediction of genetic value

$$y = g + e, g = Zu$$

$$V(u) = I\sigma_u^2, V(Zu) = ZZ'\sigma_u^2,$$

$$\hat{u} = E(u | y) = b'y = \sigma_u^2 Z'V^{-1} y$$

$$\hat{g} = Z \hat{u} = \sigma_u^2 ZZ'V^{-1} y = \sigma_g^2 GV^{-1} y$$

29

## Prediction of genetic value

$$y = g + e, g = Zu$$

If  $y$  and  $u$  are multivariate normal

$$E(u | y) = b'y = \sigma_u^2 Z'V^{-1} y$$

The SNP effects are unlikely to be normally distributed with equal variance

30

## Prediction of genetic value

### Best prediction

$$\begin{aligned} \hat{u} &= E(u | y) \\ &= \int u P(u | y) du \end{aligned}$$

Bayes theorem

$$P(u | y) = P(y | u) P(u) / P(\text{data})$$

↑                    ↓  
Likelihood      prior

31

## Prediction of genetic value

### Bayesian estimation

$$E(u | y) = \int u P(y | u) P(u) / P(y) du$$

Distribution of SNP effects

Normal	→ BLUP
t-distribution	→ Bayes A
Mixture	→ Bayes B (Meuwissen et al 2001)

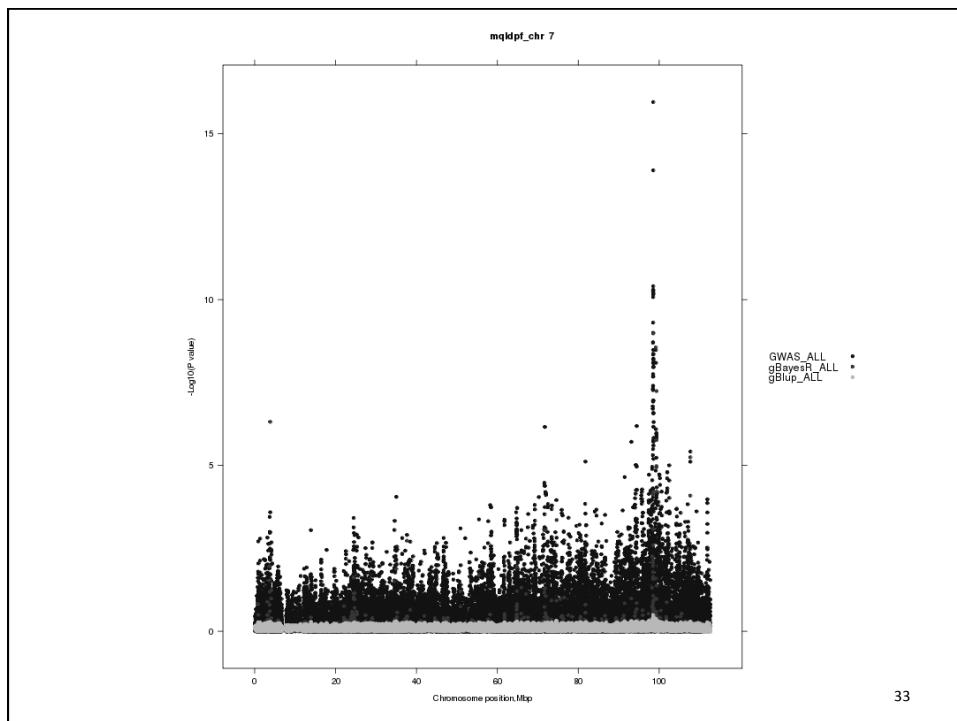
Mixture of N → Bayes R (Erbe et al 2012)

$u \sim N(0, \sigma_i^2)$  with probability  $\pi_i$

$$\sigma_i^2 = \{0, 0.0001, 0.001, 0.01\} \sigma_g^2$$

Accuracy is greatest if assumed distribution matches real distribution.

32



## Prediction of genetic value

Other methods of prediction

Estimate effect of each SNP one at a time and add  
 $\hat{g} = Z \hat{u}$   
 $\hat{u}$  estimated from single SNP regression

Biased  $E(g | \hat{g}) \neq g$   
 Less accurate because ignores LD between SNPs  
 and treats  $u$  as fixed effects

# Prediction of genetic value

## Real data

4500 bulls and 12000 cows (Holstein and Jersey)

600,000 SNPs genotyped

Train using bulls born < 2005

Test using bulls born >= 2005

Correlation of EBV and daughter average

	Protein	Stature	Milk	Fat%
BLUP	0.66	0.52	0.65	0.72
Bayes R	0.66	0.54	0.68	0.82

35

## Genetic architecture



Proportion of SNPs from distribution with variance

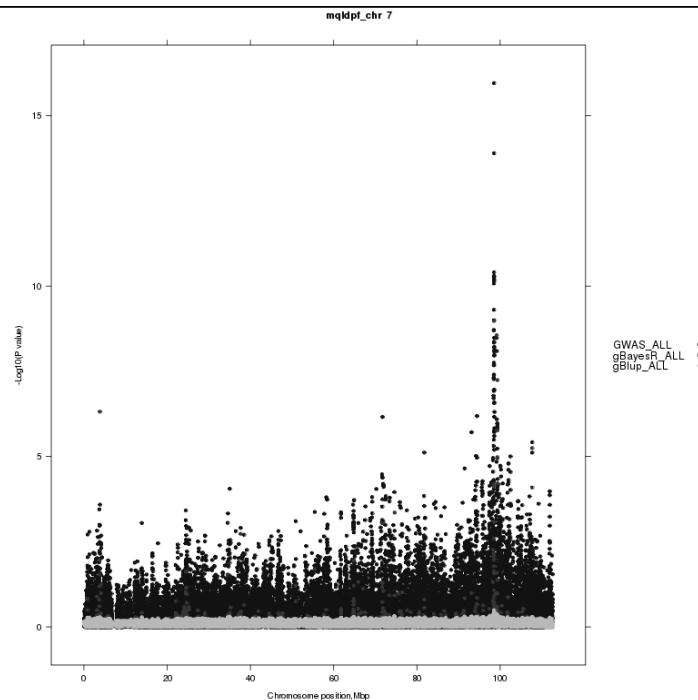
Trait	0.01%	0.1%	1%	polygenic (%)
RFI	7498	296	6	11
LDPF	1419	254	36	27
Mean	4029	271	19	25

36

## Integration of prediction and mapping of causal variants

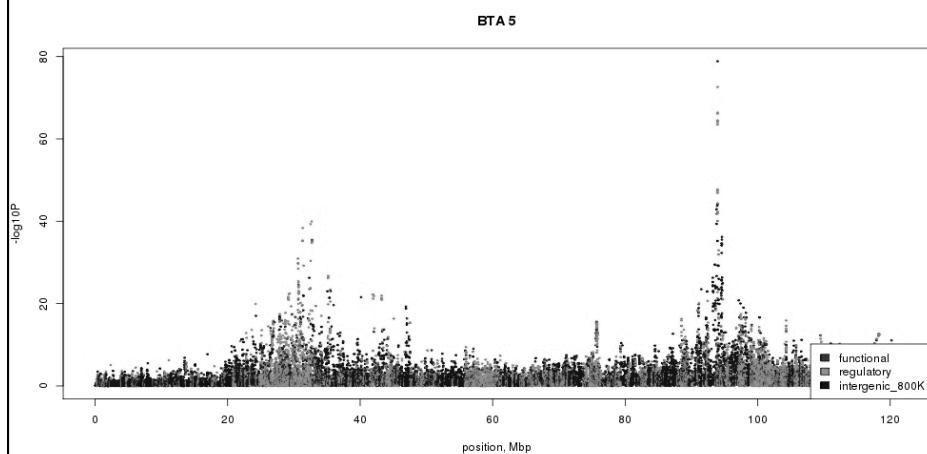
Same Bayesian models as used for prediction  
can be used for mapping causal variants of  
complex traits

37



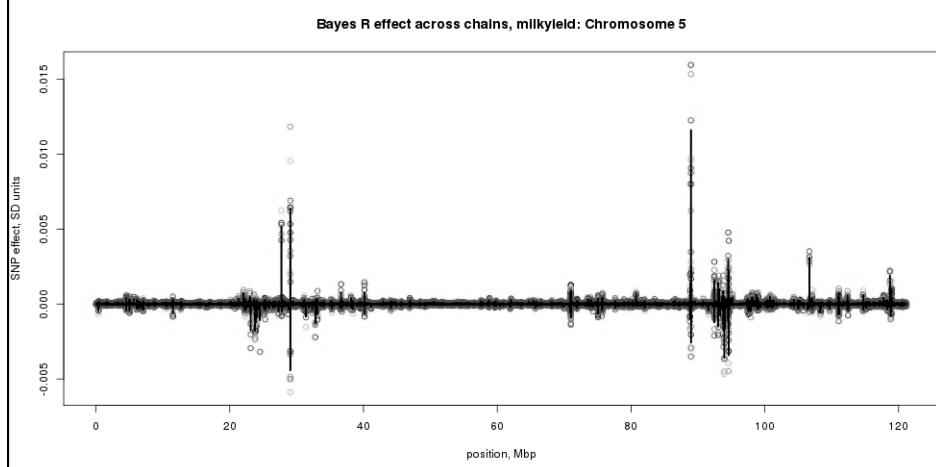
38

## Mapping QTL – Milk on BTA5



39

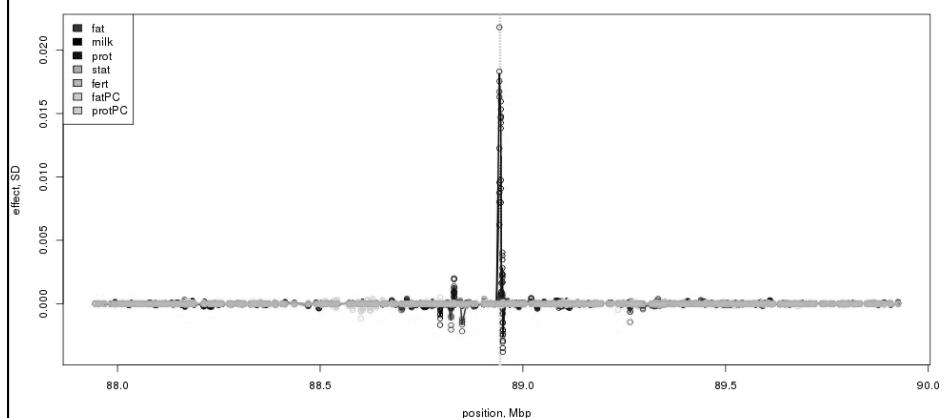
## Mapping QTL – Milk on BTA5



40

## Mapping QTL – Milk on BTA5

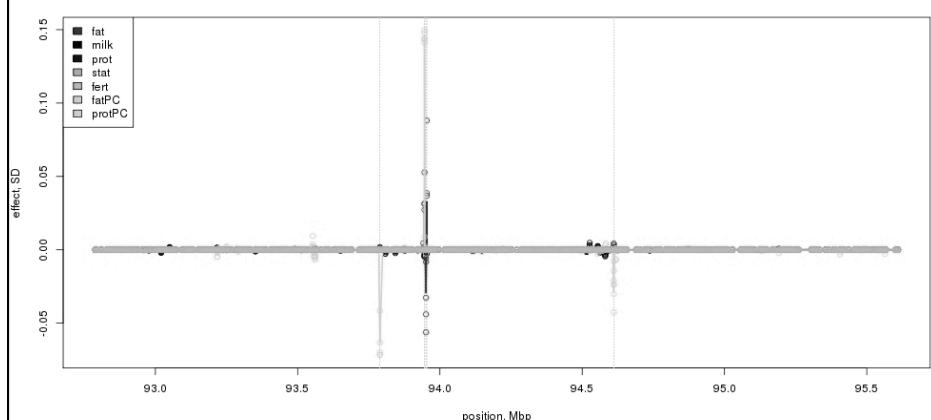
region 10 : chromosome 5 , 1.98 Mb, centered on 88.94 ( 13 hit/s)



41

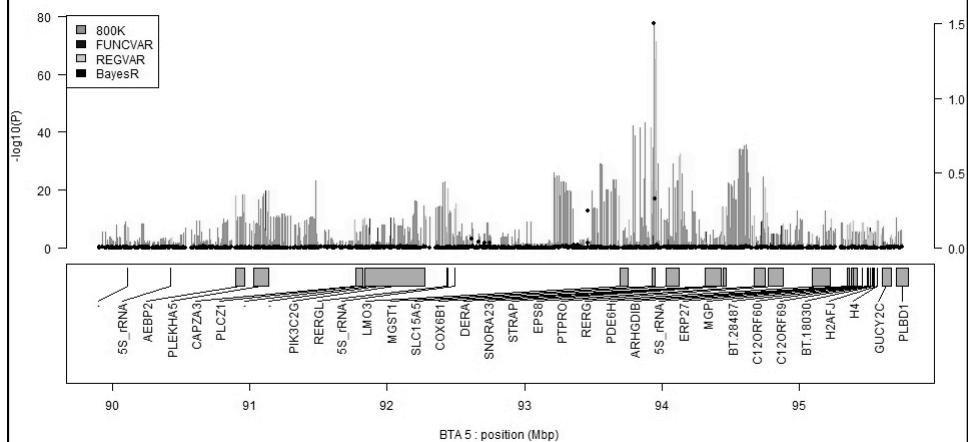
## Mapping QTL - Milk on BTA5

region 11 : chromosome 5 , 2.82 Mb, centered on 94.2 ( 26 hit/s)



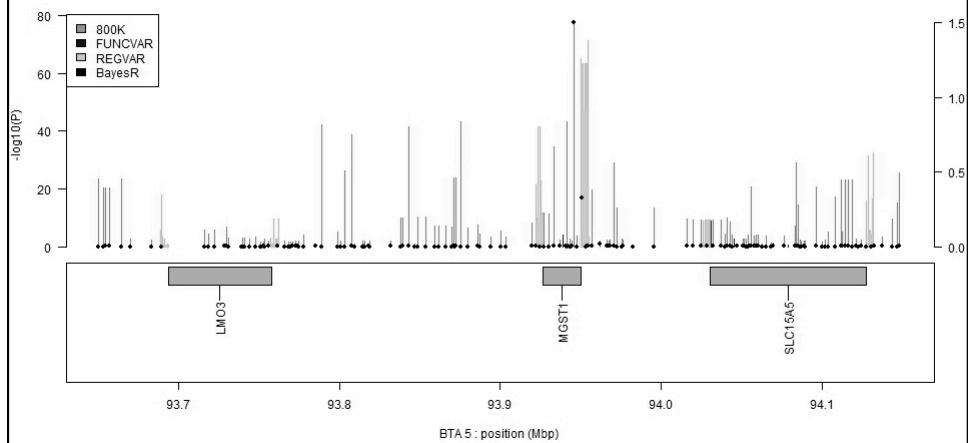
42

## Mapping QTL – Milk on BTA5



43

## Mapping QTL – Milk on BTA5



44

## Prediction of genetic value Summary

Best prediction is  $\hat{g} = E(g | y)$

Genetic values treated as random effects

$$Eg \quad g \sim N(0, G\sigma_g^2)$$

Equivalent model to predict SNP effects  $u$

$E(u | y)$  depends on prior distribution of  $u$

→ Bayesian models

$\hat{g} = Z \hat{u}$  gives higher accuracy than assuming

$$g \sim N(0, G\sigma_g^2)$$

Bayesian models integrate prediction and mapping of causal variants

45

## Key concepts

- Prediction of phenotypic values is limited by heritability
- Accuracy of prediction depends on
  - how well marker effects are estimated (sample size)
  - how well marker effects are correlated with causal variants (LD)
- Estimation of marker effects and prediction in the same data leads to (severe) bias
  - winner's curse; over-fitting
- Variance explained by a SNP-based predictor is not the same as the variance explained by those SNPs
- Marker data captures both between and within family genetic variation
- Best prediction methods take genetic values as random effects

46