

# **SISG Brisbane Module 10: Statistical & Quantitative Genetics of Disease**

## ***Lecture 3*** ***Single locus model of disease risk*** ***Naomi Wray***

# ***Aims of Lecture 3***

## Theory

- Single locus disease model
- Power calculations

# *Single locus disease model*

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population;

p = risk allele frequency;

	P(G)
aa	$(1-p)^2$
Aa	$2p(1-p)$
AA	$p^2$

# Single locus disease model

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population;

p = risk allele frequency;

$f_0$  = baseline risk for homozygote non-risk allele – UNKNOWN

R = relative risk for heterozygote; assume risk is multiplicative (on this scale)

	P(G)	P(D G)
aa	$(1-p)^2$	$f_0$
Aa	$2p(1-p)$	$f_0R$
AA	$p^2$	$f_0R^2$

# Single locus disease model

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population;

p = risk allele frequency;

f<sub>0</sub> = baseline risk for homozygote non-risk allele – UNKNOWN

R = relative risk for heterozygote; assume risk is multiplicative (on this scale)

	P(G)	P(D G)	P(D) =P(D G)p(G)
aa	(1-p) <sup>2</sup>	f <sub>0</sub>	(1-p) <sup>2</sup> f <sub>0</sub>
Aa	2p(1-p)	f <sub>0</sub> R	2p(1-p) f <sub>0</sub> R
AA	p <sup>2</sup>	f <sub>0</sub> R <sup>2</sup>	p <sup>2</sup> f <sub>0</sub> R <sup>2</sup>
			Sum= K

$$P(\text{Disease})=K = f_0(1-p)^2 + f_0R2p(1-p) + f_0R^2p^2 = f_0(1+p(R-1))^2$$

$$f_0=K/(1+p(R-1))^2$$

# Single locus disease model

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population;

p = risk allele frequency;

f<sub>0</sub> = baseline risk for homozygote non-risk allele – UNKNOWN

R = relative risk for heterozygote; assume risk is multiplicative (on this scale)

	P(G)	P(D G)	P(D) =P(D G)p(G)	P(G D) =P(G)/P(D)
aa	(1-p) <sup>2</sup>	f <sub>0</sub>	(1-p) <sup>2</sup> f <sub>0</sub>	(1-p) <sup>2</sup> f <sub>0</sub> /K
Aa	2p(1-p)	f <sub>0</sub> R	2p(1-p) f <sub>0</sub> R	2p(1-p) f <sub>0</sub> R/K
AA	p <sup>2</sup>	f <sub>0</sub> R <sup>2</sup>	p <sup>2</sup> f <sub>0</sub> R <sup>2</sup>	p <sup>2</sup> f <sub>0</sub> R <sup>2</sup> /K
			Sum= K	

$$P(\text{Disease})=K = f_0(1-p)^2 + f_0R2p(1-p) + f_0R^2p^2 = f_0(1+p(R-1))^2$$

$$f_0=K/(1+p(R-1))^2$$

# Practical

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population = 0.01;

p = risk allele frequency = 0.2;

f<sub>0</sub> = baseline risk for homozygote non-risk allele – UNKNOWN

R = relative risk for heterozygote; assume risk is multiplicative (on this scale) = 1.2

	P(G)	P(D G)	P(D) =P(D G)p(G)	P(G D) =P(G)/P(D)
aa				
Aa				
AA				

$$P(\text{Disease})=K = f_0(1-p)^2 + f_0R^2p(1-p) + f_0R^2p^2 = f_0(1+p(R-1))^2$$

$$f_0=K/(1+p(R-1))^2$$





***Using the single locus disease model to calculate power in an association study***

# What is power?

When we set up a statistical test

- The null hypothesis is EITHER
  - true
  - false
- With the data available we EITHER
  - reject the null hypothesis
  - fail to reject the null hypothesis

	<b>Null hypothesis is true</b>	<b>Null hypothesis is false</b>
Reject the null hypothesis	Type I error <b>False positive</b>	Correct Outcome <b>True positive</b>
Fail to reject the null hypothesis	Correct Outcome <b>True negative</b>	Type II error <b>False Negative</b>

Power = probability of rejecting the null hypothesis when the null hypothesis is false

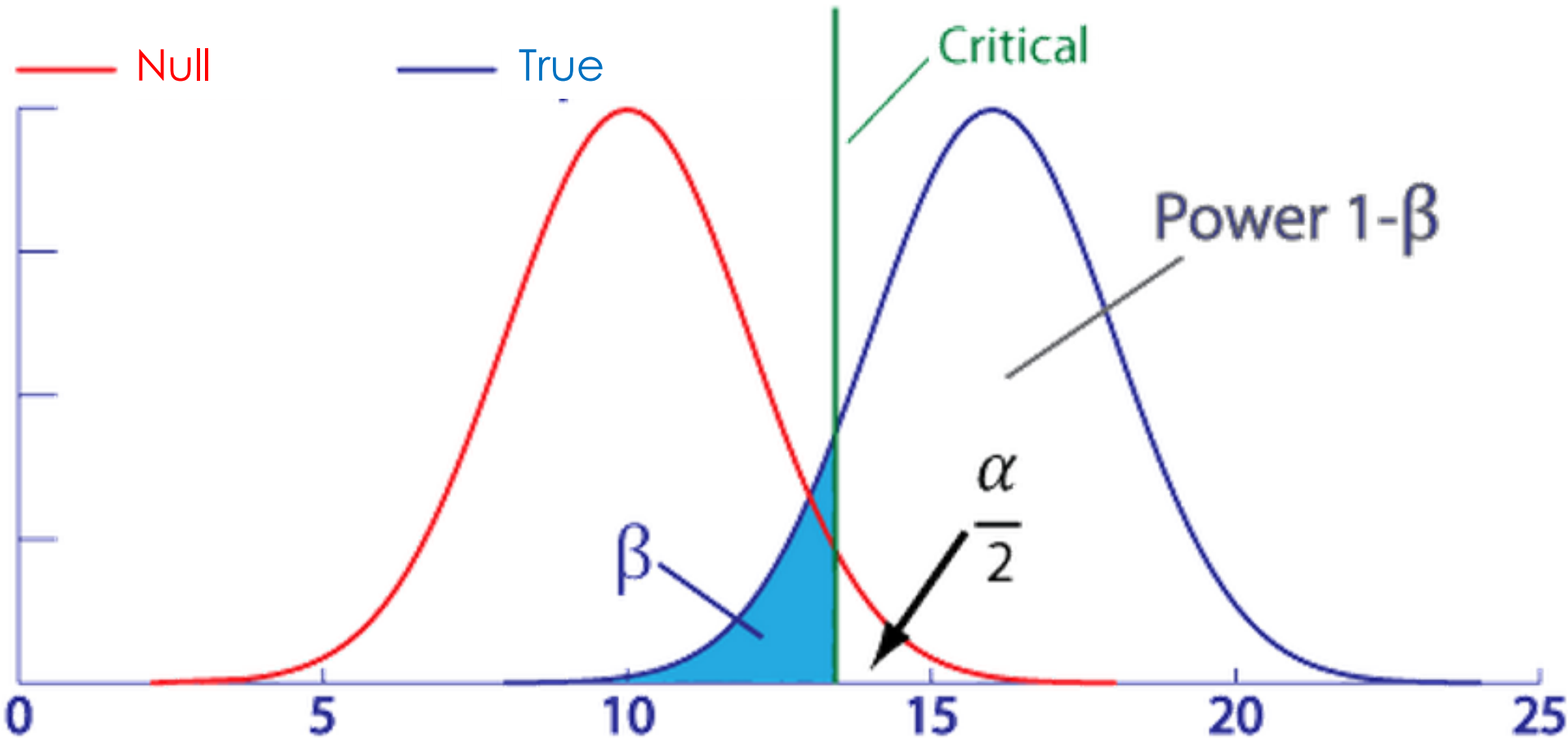
= 1 - probability of failing to reject the null hypothesis when the null hypothesis is false

= 1 - probability(Type II error)

Power depends on statistical test, effect size to be detected, sample size, acceptable level of Type I error

Non-centrality parameter depends on statistical test, effect size to be detected, sample size

# Power

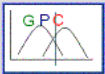


$\beta$  = probability of rejecting the null hypothesis when the alternative hypothesis is true

$\alpha$  = probability of rejecting the null hypothesis when the null hypothesis is true

Variance about mean values depends on sample size

# 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)

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

### Quantitative Case-Control

Total QTL variance :  (0 - 1)  
Dominance : additive QTL effects :  (0 - 1)  
QTL increaser allele frequency :  (0 - 1)  
Marker M1 allele frequency :  (0 - 1)  
Linkage disequilibrium (D-prime) :  (0 - 1)  
Number of cases :  (> 0)  
Case lower threshold :   
Case upper threshold :   
Control:case ratio :  (> 0)  
Controls lower threshold :   
Controls upper threshold :   
User-defined type I error rate :  (0.00000001 - 0.5)  
User-defined power: determine N :  (0 - 1)  
(1 - type II error rate)

## Genetic Power Calculator

### Case - control for discrete traits

High risk allele frequency (A) :  (0 - 1)  
Prevalence :  (0.0001 - 0.9999)  
Genotype relative risk Aa :  (>1)  
Genotype relative risk AA :  (>1)  
D-prime :  (0 - 1)  
Marker allele frequency (B) :  (0 - 1)  
Number of cases :  (0 - 10000000)  
Control : case ratio :  (> 0)  
( 1 = equal number of cases and controls)  
 Unselected controls? (\* see below)  
User-defined type I error rate :  (0.00000001 - 0.5)  
User-defined power: determine N :  (0 - 1)  
(1 - type II error rate)

Created by [Shaun Purcell](#) 24.Oct.2008

# Genetic Power Calculator

## Case - control for discrete traits

High risk allele frequency (A) :  ( 0 - 1 )  
Prevalence :  ( 0.0001 - 0.9999 )  
Genotype relative risk Aa :  ( >1 )  
Genotype relative risk AA :  ( >1 )  
D-prime :  ( 0 - 1 )  
Marker allele frequency (B) :  ( 0 - 1 )  
Number of cases :  ( 0 - 10000000 )  
Control : case ratio :  ( >0 )  
( 1 = equal number of cases and controls)

Unselected controls? (\* see below)

User-defined type I error rate :  ( 0.00000001 - 0.5 )

User-defined power: determine N :  ( 0 - 1 )

( 1 - type II error rate )

## Case-control statistics: allelic 1 df test (B versus b)

Sample NCP = 28.59

Alpha	Power	N cases for 80% power
0.1	0.9999	1081
0.05	0.9996	1372
0.01	0.9972	2042
0.001	0.9802	2985
5e-08	0.4586	6924

# Power of a case-control study

Power of a disease trait

$p$  = frequency of risk allele in population

$p_{case}$  = frequency of risk allele in cases

$p_{cont}$  = frequency of risk allele in controls

$v$  = proportion of a sample of  $N$  that are cases

$\bar{p}$  = mean allele frequency across cases and controls  
=  $v p_{case} + (1-v) p_{control}$

--

# Power of a case-control study

Power of a disease trait

$p$  = frequency of risk allele in population

$p_{case}$  = frequency of risk allele in cases

$p_{cont}$  = frequency of risk allele in controls

$v$  = proportion of a sample of  $N$  that are cases

= mean allele frequency across cases and controls

$\bar{p}$  =  $v p_{case} + (1-v) p_{control}$

Z-Test statistic of association = test of difference of two proportions =

$$\frac{p_{case} - p_{cont}}{s.e. (pooled sample p)} = \frac{p_{case} - p_{cont}}{s.e. (\bar{p})}$$

$$\chi^2 \text{ non-centrality parameter} = NCP_{01} = \frac{(p_{case} - p_{cont})^2}{var(\bar{p})}$$

$$var(\bar{p}) = \frac{1}{2} \bar{p}(1 - \bar{p}) \left( \frac{1}{Nv} + \frac{1}{N(1 - v)} \right)$$

# Allele Frequency in Cases

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population

	P(G)	P(D G)	P(D) =P(D G)p(G)	P(G D) =P(G)/P(D)
aa	$(1-p)^2$	$f_0$	$(1-p)^2 f_0$	$(1-p)^2 f_0/K$
Aa	$2p(1-p)$	$f_0R$	$2p(1-p) f_0R$	$2p(1-p) f_0R/K$
AA	$p^2$	$f_0R^2$	$p^2 f_0R^2$	$p^2 f_0R^2/K$
			Sum= K	

$$P(\text{Disease})=K = f_0(1-p)^2 + f_0R2p(1-p) + f_0R^2p^2 = f_0(1+p(R-1))^2$$

$$f_0=K/(1+p(R-1))^2$$

$$\begin{aligned} p_{\text{case}} &= \frac{1}{2} P(\text{Aa}|\text{D})+P(\text{AA}|\text{D}) \quad \text{Allele frequency in cases} \\ &= f_0pR((1-p) + pR)/K = \frac{pR}{(1+p(R-1))} \end{aligned}$$

Find allele frequency in controls in the same way

$$p_{\text{cont}} = \frac{p}{1-K} \left( 1 - \frac{KR}{(1+p(R-1))} \right)$$



# Allele Frequency in Controls

Single locus disease model:

G = genotype; D=disease; K = overall disease risk in population

	P(G)	P(D' G)	P(D') =P(D' G)p(G)	P(G D') =P(G)/P(D')
aa	$(1-p)^2$	$(1-f_0)$	$(1-p)^2 (1-f_0)$	$(1-p)^2 (1-f_0)/(1-K)$
Aa	$2p(1-p)$	$(1-f_0R)$	$2p(1-p) (1-f_0R)$	$2p(1-p) (1-f_0R)/(1-K)$
AA	$p^2$	$(1-f_0R^2)$	$p^2 (1-f_0R^2)$	$p^2 (1-f_0R^2)/(1-K)$
			Sum= 1-K	

$$f_0 = K / (1 + p(R-1))^2$$

$$p_{\text{control}} = \frac{1}{2} P(\text{Aa}|D') + P(\text{AA}|D') \quad \text{Allele frequency in controls}$$

$$= \frac{p}{1-K} \left( 1 - \frac{KR}{(1+p(R-1))} \right)$$

# Power of a case-control study

$$NCP_{01} = \frac{(p_{case} - p_{cont})^2}{var(\bar{p})}$$

$\alpha$  = significance level - acceptable level of type I error

$t = \Phi^{-1}\left(\frac{\alpha}{2}\right)$  Normal distribution threshold above which null hypothesis will be rejected

$$Power = \Phi(\sqrt{NCP_{01}} + t)$$

$N=10000, v=0.5, p=0.2, R=1.2, K=0.01, \alpha=5e-8, K=0.01, power = 0.46$

Agrees with the genetic power calculator

## **Yang et al (2009) Comparing Apples and Oranges: Equating the Power of Case-Control and Quantitative Trait Association Studies. Genetic Epidemiology**

Research Question in 2009:

We had GWAS success with height but not with disease.

Was this a function of power?

For the same sample size what is the connection between power for a quantitative trait vs case-control?

Answer:

- “So a planned meta-analysis for height on 120,000 individuals has power equivalent to a CC study on 33,100 schizophrenia cases and 33,100 controls, a size not yet achievable for this disease.”

# Power of a case-control association study expressed in terms of variance explained by the locus

$$\chi^2 \text{ non-centrality parameter} = NCP_{01} = \frac{(p_{case} - p_{cont})^2}{var(\bar{p})}$$

$$NCP_{01} = \frac{2\bar{p}(1 - \bar{p})(R - 1)^2 v(1 - v)N}{(1 - K)^2 (1 + p(R - 1))^2}$$

If  $R$  is small then  $(1 + p(R - 1))^2 \approx 1$  e.g.,  $p = 0.2$ ,  $R = 1.2$ ,  $(1 + p(R - 1))^2 = 1.08$

$$\text{Variance explained by a locus} = h_{L[j]}^2 \approx \frac{2p(1 - p)(R - 1)^2}{i^2}$$

$$NCP_{01} \approx \frac{h_{L[j]}^2 i^2 v(1 - v)N}{(1 - K)^2}$$

# Approximate variance explained by a locus

Regression of disease on  $j$ th SNP,  $x_{[j]} = 0, 1, 2$

$$y_{01} = K + b_{01}x_{[j]} + \varepsilon$$

When  $x_{[j]}=0$   $\hat{y}_{01} = K$  = P(Disease | Genotype = aa)

When  $x_{[j]}=1$   $\hat{y}_{01} = K + b_{01}$  = P(Disease | Genotype = Aa)

Relative Risk = R = P(Disease | Genotype = Aa) / P(Disease | Genotype = aa)

$$= (K + b_{01}) / K \quad \text{so} \quad b_{01} = K(R - 1)$$

Variance attributable to the locus on the disease scale

$$\sigma_{A_{01}[j]}^2 = h_{01[j]}^2 K(1 - K) = b_{01}^2 \text{var}(x) = 2p(1 - p)b_{01}^2$$

$$h_{01[j]}^2 = 2p(1 - p)b_{01}^2 / K(1 - K)$$

$$h_{L[j]}^2 = \frac{(1 - K)h_{01[j]}^2}{i^2 K} = \frac{2p(1 - p)b_{01}^2}{i^2 K^2} = \frac{2p(1 - p)(R - 1)^2}{i^2}$$

Assumes a population sample not a case control sample



# ***Power of a association study of a quantitative trait***

$$\chi^2 \text{ non-centrality parameter} = NCP_{QT} = \frac{N_{QT} h_{L[i]}^2}{1 - h_{L[i]}^2}$$

***When the variance explained is the same in c-c and for quantitative trait***

$$NCP_{01} \approx \frac{h_{L[j]}^2 i^2 v(1 - v) N_{01}}{(1 - K)^2}$$

$$\frac{NCP_{01}}{NCP_{QT}} \approx \frac{i^2 v(1 - v) N_{01}}{(1 - K)^2 N_{QT}}$$

# Practical

- a) Code for slides 3-6. Done already
- b) Power in case-control study design
  - i) Compare to GPC
  - ii) Compare power for screened to unscreened controls
  - iii) Compare the impact on power of screening controls for # schizophrenia  $K = 0.01$  and Major depression  $K = 0.15$
  - iv) For which disorders is screening of controls most recommended
- c) Power Graphs in case-control study design
  - First run code through to see graphs, then look at code. Makes graph with 3 lines based on RAF for one disease risk
    - i) You make graph with 3 lines based on disease risk  $K=0.001, 0.01, 0.1$
  - # Q: For a given GRR is the power bigger or smaller as disease prevalence increases
  - # Q: Why does this make sense?
- d) Just run code