2017 SISG Brisbane Module 10: Statistical & Quantitative Genetics of Disease

Lecture 5 Polygenic models of disease risk Naomi Wray

Aims of Lecture 5

- To consider polygenic models of genetic risk
- To demonstrate that many polygenic models are consistent with empirical data and that they can be considered equivalent
- To understand the conclusion that the liability threshold model is the model of choice
- To understand the criticisms and controversy of the liability threshold model
- Introduction to polygenic risk scores

Genetic models of disease

Mendelian disease:

- Individuals that possess the mutation get the disease.
- Dominant e.g Huntington's or recessive e.g. Cystic fibrosis

Mendelian disease with variable penetrance.

- Only those with the mutation get the disease
- Not everyone with the mutation gets the disease.
- E.g. C9orf72 in Motor Neurone Disease

Compound heterozygote disease.

• Like recessive Mendelian but individuals carry two different rare mutations in the same gene.

Two-hit diseases

• Hypothesized, but examples?

Oligogenic diseases – caused by presence of several genetic risk variants Polygenic diseases – caused by multiple genetic risk variants Multifactorial diseases- caused by multiple genetic risk variants and other risk factors

Common complex genetic diseases are likely to be polygenic multifactorial

Evidence:

Many risk variants of small effect identified

Implications:

- We all carry risk alleles
- Each affected person may carry a unique portfolio
- Polygenic model can accommodate some people having few loci of larger effect and others having many loci of small effect
- The more loci involved, to be consistent with low prevalence, the probability of disease has to increase steeply with the number of loci.
- The more loci involved, the more likely they have a pleiotropic effect, which would be consistent with them being common in the population
- The more loci involved implies that we are highly robust to perturbations – but this breaks down when the burden of risk factors become too great

Modeling polygenic genetic risk

- "Easiest" to understand by thinking of individual risk loci and how they act together to cause disease
 - The frequency of the risk alleles
 - Drawn from a distribution
 - All the same
 - The effect size of the risk alleles
 - Drawn from a distribution
 - All the same relative risk associated
 - Interaction between risk loci
 - Complex
 - All act in the same way



Start simple

- Assume all risk alleles have the same frequency
- Assume all risk alleles have the same effect size
- Then risk is a reflection of the count of the number of risk alleles
- What is the shape of the relationship between count of risk alleles and probability of disease?

Count risk alleles

 SNP 1: Aa
 1

 SNP 2: Aa
 1

 SNP 3: AA
 2

 SNP 4: aa
 0

 SNP 5: aa
 0

4 risk alleles =s

Basic Model

0.1

100

p = freq of risk allele 1-p = freq of non-risk allele Assume Hardy- Weinberg equilibrium in the population Genotype frequencies $P(bb) = (1-p)^2$ P(Bb) = 2p(1-p)

 $P(BB) = p^2$

Relative risk associated with one risk allele R

n loci

Theoretical minimum number of risk loci : 00Theoretical maximum number of risk loci possible: 2n200

Mean number of risk loci: 2np20Variance in number of risk loci: 2np(1-p)18Range in number of loci expected 2np +/- $(3.5)\sqrt{2np(1-p)}$ 5 - 36

Visualising common complex genetic diseases Polygenic genetic architecture

- Imagine a disorder underpinned by
 - 100 loci : 2 alleles at each locus
 - Each risk allele has frequency 0.1



- 0 risk alleles = yellow
- 1 risk allele = light blue
- 2 risk alleles = dark blue

Average person a person carries 2 alleles * 100 loci *0.1 = 20 risk alleles

Everybody carries some risk alleles

Range in population \sim 5-36 (mean +/- 3.5 sd)

Polygenic burden : top 1% carry > 33 risk alleles

Visualising variation between individuals for common complex genetic diseases



Not all affected individuals carry the risk allele at any particular locus Unaffected individuals carry multiple risk loci Consequences of risk alleles depend on the genetic and environmental background

How to combine risk loci to explain disease

Additive on disease scale

Multiplicative on disease scale

Constrained multiplicative on disease scale

Multiplicative Odds on disease scale

Liability threshold model

Basic single locus risk model

Single locus disease model:

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

risk alleles

eles		P(G)	P(D G)	P(D)=P(D G)p(G)	P(G D)=P(G)/P(D)
0	aa	$(1-p)^2$	f ₀	$(1-p)^2 f_0$	$(1-p)^2 f_0/K$
1	Aa	2p(1-p)	$f_0 R_{Bb}$	$2p(1-p) f_0 R_{Bb}$	2p(1-p) f ₀ R _{Bb} /K
2	AA	p ²	$f_0 R_{BB}$	$p^2 f_0 R_{BB}$	$p^2 f_0 R_{BB}/K$
				Sum= K	

 $P(Disease) = K = f_0(1-p)^2 + f_0 R_{Bb} 2p(1-p) + f_0 R_{BB} p^2 = f_0(1+p(R-1))^2$

$$= f_0((1-p)^2 + R_{Bb}2p(1-p) + R_{BB} p^2)$$

$$f_0 = K/((1-p)^2 + R_{Bb}2p(1-p) + R_{BB}p^2)$$

if $R_{Bb} = R$; $R_{BB} = R^2$ multiplicative on disease scale

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

 $R_{BB} = 2^* R_{Bb}$ additive on the disease scale

Two-locus risk model

Single locus disease model:

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

P(G)

P(D)

# risk	$1^{st} locus/2^{nd} locus$	Frequency	Risk Additive	Risk
alleles				Multiplicative
0	0/0	$(1-p)^4$	f_0	f_0
1	1/0 or 0/1	$2p(1-p)(1-p)^2 + 2p(1-p)(1-p)^2$	f ₀ R	f ₀ R
2	2/0 or 0/2 or 1/1	$p^2 + p^2 + 4p^2(1-p)^2$	$2f_0R$	$f_0 R^2$
3	2/1 or 1/2	$p^{2}2p(1-p) + p^{2}2p(1-p)$	$3f_0R$	$f_0 R^3$
4	2/2	p ⁴	$4f_0R$	$f_0 R^4$

Sum P(D) = K

 f_0 = probability of disease with no risk alleles. This baseline probability differs between the models

Additive on the disease scale

Probability of disease increases additively/linearly with the number of loci (x) carried.

 $P(D | x = s) = b^*R^*s$

Constraint

$$\sum_{x=0}^{2n} P(D|x)P(x) = K$$

$$E(P(D | x)) = E(b^*R^*x) = b^*R^*E(x) = b^*R^*2np = K$$

So b = K/2npR



Looking at the additive model

Base

- N = 1e5 # number of families 100,000
- n = 100 # number of loci
- R = 1.1 # relative risk of each risk allele
- p = 0.2 # allele frequency of each risk allele
- K = 0.01 # probability of disease

Follow up:

Base, R=1.5, p=0.5, K =0.1

Look at maximum probability of disease and consider whether this model will generate an increased risk in relatives

SWITCH TO R

Histogram of # risk alleles Histogram of probability of dise



onship between # alleles & prob

Histogram of # risk alleles



Histogram of # risk alleles Histogram of probability of dise





Histogram of # risk alleles



indA = # risk alleles

110

130

Additive model

- Mathematically tractable
- To achieve additivity of risk loci and correct disease prevalence, does not give high probability of disease with large number of risk loci
- Not consistent with high heritability
- Not consistent with observed risks to relatives

- Can "fudge" the additive model by saying
 - P(D | x < n1) = 0
 - P(D|n1 < x < n2) = additive with x
 - P(D | x> n2) = 1

Is non-linear with x Not mathematically tractable

Multiplicative on the disease scale

Probability of disease increases multiplicatively with the number of risk loci (x)

 $P(D | x = s) = f_0 R^s$ When s =0, $P(D | x = 0) = f_0$ Multiplicative on the risk scale

Constraint

$$\sum_{x=0}^{2n} P(D|x)P(x) = K$$

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

Additive on the log risk scale

 $Log(P(D | x=s)) = s log(f_0R)$

Looking at the multiplicative model

Base

- N = 1e5 # number of families
- n = 100 # number of loci
- R = 1.1 # relative risk of each risk allele
- p = 0.2 # allele frequency of each risk allele
- K = 0.01 # probability of disease

```
Follow up:
Base, K=0.1
Base K = 0.1, R =1.2
Look at maximum probability of disease and consider whether this
model will generate an increased risk in relatives
```

Add fix

SWITCH to R

K=0.1, p=0.2, R=1.1

Histogram of # risk alleles Histogram of probability of dise









K=0.1, p=0.2, R=1.2

Histogram of # risk alleles Histogram of probability of dise 8e+04 25000 Frequency Frequency 4e+04 10000 00+90 0 20 40 60 5 10 15 0 indA = # risk alleles indR = probability of disease given # risk a

onship between # alleles & prob

indR = probability of disease Not diseased Diseased 12000 15 8000 Frequency 9 4000 S 0 0 70 30 40 50 60 70 20 30 40 50 60 20

 $ind \Delta = # risk alleles$

indA = # risk alleles

Histogram of # risk alleles

Multiplicative model

- Mathematically tractable
- High probability of disease with large number of risk loci so consistent with high heritability and can be consistent with observed risks to relatives

BUT

• Probability of disease for an individual can be > 1

IF constrain so that max probability of disease is 1 THEN

- E(P(D | x)) is no longer K
- Need to fudge to retain this property
- Loses mathematical tractability

Epidemiology risk model

Odds(Disease) = P(Disease)/(1-P(Disease))

Odds(Disease | x = s) = Odds(Disease | x = 0) $\gamma^{x} = C\gamma^{x}$

s = number of risk loci carried by an individuals

 γ = odds ratio for each risk locus

 $P(Disease | x = s) = C\gamma^{s}/(1-C\gamma^{s})$

Good: probability of disease does not exceed 1 Bad: mathematically intractable

Janssen et al (2006) Predictive testing for complex diseases using multiple genes: Fact or fiction? Genet Med 8 395 Lu & Elston (2008) Using the optimal ROC to design a predictive test, exemplified with Type 2 Diabetes AJHG 82

Epidemiology risk modelling

- R = risk = probability of disease
- $\log R = y \sim N(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$
- R ~ LogNormal($\boldsymbol{\mu}, \boldsymbol{\sigma}^2$) = LN($\boldsymbol{\mu}, \boldsymbol{\sigma}^2$)
- μ is arbitrary but Pharaoh set as $\mu = -\sigma^2/2$, but can also be calculated from disease prevalence K

$$\sigma^2 = \log(\lambda_{MZ}) = 2\log(\lambda_{sib})$$

$$\mu = \log K - \sigma^2/2$$

Epidemiology risk model

$$E[R] = K = \int e^{\mu + \sigma x} \phi(x) dx = \dots = e^{\mu + \sigma^2/2}$$

Two relatives, with risk of disease R_1R_2

$$R_1 = e^{\mu + \sigma z_1}$$

$$R_2 = e^{\mu + \rho \sigma z_1 + \sqrt{(1 - \rho^2)\sigma z_2}}$$

Probability that both are affected R_1R_2

$$E[\mathbf{R}_1\mathbf{R}_2] = \int R_1 \phi(z_1) R_2 \phi(z_2) dz_1 \ dz_2 = \dots = e^{2\mu + \sigma^2(1+\rho)}$$

Recurrence risk =
$$\lambda_{relative} = \frac{E[R_1R_2]}{K^2} = e^{\rho\sigma^2}$$

$$\lambda_{MZ} = e^{\sigma^2}$$

$$\lambda_{sib} = e^{0.5\sigma^2}$$

$$\sigma^2 = \log(\lambda_{MZ}) = 2\log(\lambda_{sib})$$

$$\mu = log K - \sigma^2/2$$





Liability threshold model

Doesn't directly parameterise in terms of number of risk loci

Only parameterises in terms of

- prevalence of disease and heritability of liability

 OR

- prevalence of disease and risk to relatives

i.e.

 In terms of total variance explained which could cover a range of genetic architectures – so can be indirectly parameterised in terms of number of risk loci

Variance explained by a locus depends on frequency (p) and effect size(a) : 2p(1-p)a²

Variance explained is the same for p=0.1, a=0.1 as for p=0.5, a=0.06

• BUT is the liability threshold model realistic?

Controversy – the abrupt threshold is not biological



"Contrary to the argument regarding the conservatism of the multifactorial threshold model for describing the inheritance of congenital malformations, little biological insight has resulted from the series of tautological, albeit grandiose, mathematical assumptions currently comprising the basis for this hypothesis." Melnick & Shields

The theoretical foundation of genome-wide association studies

GWAS are founded on the polygenic model of disease liability, which itself arises from an assertion of breathtaking audacity by the godfather of quantitative genetics, DS Falconer. In an attempt to demonstrate the relevance of quantitative genetics to the study of human disease, Falconer, based on work of others before him (for example, [24]), came up with a nifty solution [25]. Even though disease states are typically all-or-nothing, and even though the actual risk of disease is clearly very discontinuously distributed in the population (being dramatically higher in relatives of affected people, for example), he claimed that it was reasonable to assume that there was something called the underlying liability to the disorder that was actually continuously distributed.

Mitchell (2012) What is complex about complex disorders Genome Biol 12: 237

Edwards(1969) Familial predisposition in man, Br Med Bull

Melnick & Shields (1976) Allelic restriction: a biologic alternative to multifactorial threshold model. The Lancet Many references to the criticism in papers of the time eg Smith (1970)

Is the abrupt threshold non-biological?

- People are classed as diseased or not disease, any error in this classification, contributes of a heritability of < 1.
- Wright(1934) showed that 3 vs 4 toes in guinea pigs "cannot correspond to alternate phases of a single factor (=gene)" and used crosses to show several factors ("> 3") underly a physiological threshold
- Fraser (1976) Detailed explanation of the biology consistent with a multifactorial threshold model for cleft palate

Fraser(1976) The multifactorial/Threshold concept –uses and misuses Teratology Wright (1934) An analysis of variability in number of digits in an inbred strain of guineapig. Genetics 19 506 Wright (1934) The results of crosses between inbred strains of guinea pigs, differing in the number of. Genetics 19 537

No need to invoke abrupt threshold of phenotypic liability – instead use Probability of risk of disease under liability threshold model



"The abrupt threshold is thus conceptual rather than real and may be avoided by redefining the variance and risk function." Smith 1970



Probit model

Two parameters: disease prevalence and heritability

Probit model can be parameterised in terms of number of risk loci

Curnow (1972) The multifactorial model for the inheritance of liability to disease and its implications for risk to relatives. Biometrics Curnow & Smith (1975) Multifactorial models for familial diseases in man. J Royal Stat Soc A 138

Controversy – many models fit empirical data

"One cause of scepticism of the liability threshold model was the realization that the empirical data would also fit other models (Morton, '67; Smith, '71), such as a major gene combined with polygenic and environmental variation (Morton and MacLean, '74,a single locus with two alleles, each with incomplete penetrance (Reich et al., '72, or a heterogeneous mixture of cases determined either by a major locus with incomplete dominance and reduced penetrance or by environmental factors (Chung et al., '74, or various combinations of these (Elston and Stewart, '73; Lange and Elston, '75).

This is because the extreme tail of the distribution (which is all one can usually see when diseases are uncommon) are not good indicators of the shape of the main body of the distribution. "

Need risk to disease from relatives of different types of relatives to start to distinguish between models Not easy to collect, large sampling variances



Exchangeable models of disease

- For diseases 0.5%-2%
- High heritability
- Requires there be a large variance in risk among individuals. Consequently risk considered as a function of the number of causative alleles has to be steeply increasing.



Multiplicative model – standard model used but allows probability of disease to be >1. $P(Disease)=P(Disease | x=0)R^{x}$ Constrained multiplicative model – constrain the multiplicative model to have a maximum probability of 1

```
"Additive" model
P(Disease)=b+xR, b=-18/7 set
P(Disease)<0 to 0 and
P(Disease)>1to 1
```

Which polygenic model to use?

The liability threshold model is the model of choice because

- It is the simplest parameterization that fits the observable data
- It is mathematically tractable
- It makes least assumptions about genetic architecture

"Most models are wrong some models are useful"

From theory to data

Profile risk scoring

SNP profiling schematic



Visualising variation between individuals for common complex genetic diseases



- Not all affected individuals carry the risk allele at any particular locus •
- Unaffected individuals carry multiple risk loci
- Consequences of risk alleles depend on the genetic and environmental background

Acronyms, synonyms

- PRS
 - Polygenic risk score
- GPRS
 - Genomic profile risk score
- PGS
 - Polygenic score
- GRS
 - Genetic risk score
- Gene score
- Genetic score
- Genotypic score
- Allele score
- Profile score
- Linear predictor

- Discovery = Training
- Target = Replication = Test

Slide credit: Frank Dudbridge

Polygenic risk predictons turns SNP efects into predicted genetic risk for an individual





When SNP effects are estimated by standard GWAS there are arbitary decisions

P-value threshold of SNPs

Slide credit: Robert Maier



Purcell / ISC et al. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder Nature 2009

PRS ("profile scoring") in PLINK

¹⁴JOURNAL₀,CHILD PSYCHOLOGYandPSYCHIATRY Research Review: Polygenic methods and their application to psychiatric traits

Naomi R. Wray,¹ Sang Hong Lee,¹ Divya Mehta,¹ Anna A.E. Vinkhuyzen,¹ Frank Dudbridge,² and Christel M. Middeldorp^{3,4}

Reduce the SNPs to a set in approximate LD

- Best done using "clumping": keep the most associated SNP, remove those in LD with it, then keep the next most associated remaining SNP, etc
- Recommend LD threshold of r²<0.1
 plink --clump-p1 1 --clump-p2 1 --clump-r2 0.1 --clump-kb 500
- Estimate effect sizes of these SNPs plink --assoc

PRS ("profile scoring") in PLINK

- Create a file listing p-value thresholds to select SNPs into PRS
- Generate PRS for subjects in the target sample

plink --score --q-score-range

 Regress target phenotype against PRS glm(y~prs, family=binomial)

PRSice

- R package to simplify calculations of PRS
- http://prsice.info

Bioinformatics, 31(9), 2015, 1466-1468

PRSice: Polygenic Risk Score software

Jack Euesden*, Cathryn M. Lewis and Paul F. O'Reilly*

MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, United Kingdom

Practical – Risk Models

- 1. Additive risk model.
 - a. Run code
 - b. Change parameters
- 2. Multiplicative risk model.
 - a. Run code
 - b. Change parameters
- 3. Logistic risk model.
 - a. Run code
 - b. Change parameters
- 4. Liability threshold model
 - a. Run code
 - b. Change parameters

Practical – Polygenic risk scoring

- Demonstrate impact of overlap between Discovery & Target sample
- Learn to simulate SNP data
- Explore impact of
 - N sample size
 - M markers
 - h² true variance explained by SNPs
 - on variance explained by predictor





####