

GWAS.

9/1

* About human genome → 3 billion base pairs

* Deoxyribonucleic Acid

* Single Nucleotide Polymorphisms (SNPs)

→ A DNA variation when a single A, T, C, G differ between species or paired chromosomes

* Basic Definitions

I Locus → A specific location in genome, Could be a single DNA base to gene to bend

II Variant → A genetic difference relative to human reference genome.

III Polymorphism → A variant common in population (1% Allel, not individual-specific)

IV Allele → A genetic variant on a specific locus

V Genotype → The combination of allele in a individual

A SNP with two allele, A, a.

P = Percentage of A

$$P + q = 1. \rightarrow \text{Variant}$$

q = Percentage of a

$$P^2 + 2Pq + q^2 = 1. \rightarrow \text{Genotype}$$

↓
Could used to estimate the
Genotype derived disease incident rate

* Hardy-Weinberg Equilibrium

An ideal population, the observed genotype frequency
is identical to expected frequency.

No mutation, no assortative mating, no migration

no natural selection

→ Most population stay in HWE, $HWE \sim \chi^2$

	AA	Aa	aa	Total
observation	N_{AA}	N_{Aa}	N_{aa}	N
Expected	NP^2	Npq	Nq^2	N

$$\sum_j \left(\frac{(O_g - E_g)^2}{E_g} \right) \sim \chi^2$$

$g=AA$
 Aa
 aa

Linkage Disequilibrium

Loci are associated with each other in Non-random level.

Ex

Haplotype

A ₁ B ₁	x_{11}	$P_1 = \text{Prop of } A_1$
A ₁ B ₂	x_{12}	$P_2 = \text{Prop of } A_2$
A ₂ B ₁	x_{21}	$g_1 = \text{Prop of } B_1$
A ₂ B ₂	x_{22}	$g_2 = \text{Prop of } B_2$

$$\downarrow \quad P_1 + P_2 = g_1 + g_2 = 1$$

$$x_{11} + x_{12} = P_1.$$

$$x_{11} + x_{21} = g_1$$

Define $D = x_{11} - P_1 \cdot g_1 = P(AB) - P(A) \cdot P(B)$

↓

the deviation of expected and observed frequency

↓

Serve as Covariance

→ Next

The formula of correlation between two Variables

$$r = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} \rightarrow D$$

In Binary data

Variance is $P \cdot (1-P)$



$$r = \frac{D}{\sqrt{P_1 \cdot P_2 \cdot q_1 \cdot q_2}}$$

The correlation between
Pairs of Loci

To test LD → by Chi-square

Haplotype	Observed	Expected
A ₁ B ₁	N ₁₁	N · P ₁ · q ₁
A ₁ B ₂	N ₁₂	N · P ₁ · q ₂
A ₂ B ₁	N ₂₁	N · P ₂ · q ₁
A ₂ B ₂	N ₂₂	N · P ₂ · q ₂

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

The meaning of the "D."

$$\textcircled{1} \quad D = P(AB) - P(A) \cdot P(B) = 0$$

$P(AB) = P(A) \cdot P(B) \rightarrow$ The loci do not affect each other

$$\textcircled{2} \quad D = P(AB) - P(A) \cdot P(B) > 0$$

$P(AB) > P(A) \cdot P(B) \rightarrow$ The loci influence each other in positive

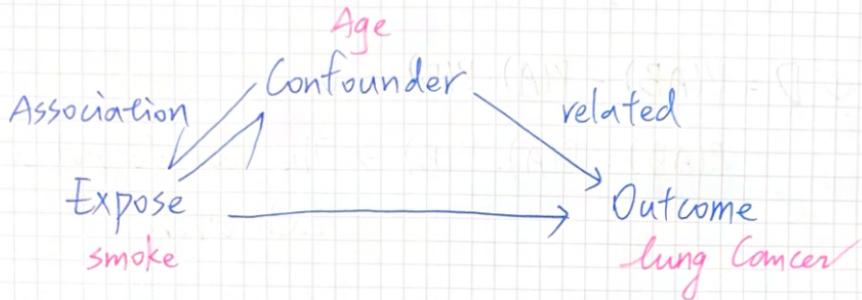
Tend to be together association

$$\textcircled{3} \quad D = P(AB) - P(A) \cdot P(B) < 0$$

$P(AB) < P(A) \cdot P(B) \rightarrow$ The loci influence each other in negative

Tend to be Separative association

Confounder



GWAS ?

A Study associated with genotype and phenotype (trait)

Type I Error \rightarrow False positive

① Avoid Post hoc subgroup analysis \rightarrow multiple test

② Statistical adjustment

I Bonferroni adjustment

II Perturbation test

$$\text{Type I Error} = 1 - (1-\alpha)^M = \alpha M$$

Why M is large, more type I Error

Type II Error False negative findings

Genotype and Environment

		Case	Control	
G+	E+	a	b	$OR_{ge} = \frac{a \cdot h}{b \cdot g}$
G+	E-	c	d	$OR_g = \frac{c \cdot h}{d \cdot g}$
G-	E+	e	f	$OR_e = \frac{e \cdot h}{f \cdot g}$
G-	E-	g	h	

Addictive interaction

$$OR_{ge} = OR_g + OR_e - 1$$

↳ minus baseline

Multiplication interaction

$$OR_{ge} = OR_g \cdot OR_e$$

Ex: Hypothesize Case is a Genotype and environment interaction disease

For Case

	G+	G-
E+	a	b
E-	c	d

$$OR = \frac{a \cdot g}{b \cdot c} \neq 1$$

For Control

	G+	G-
E+	b	f
E-	g	h

$$OR = \frac{b \cdot h}{g \cdot f} = 1$$