



Genome wide association studies

Helga Westerlind, PhD

Outline

- About GWAS/Complex diseases
- How to GWAS
- Imputation



What is a genome wide association study?

Why are we doing them?

In what context?



How do we know there is genetics involved in the disease susceptibility?

Karolinska Institutet



Family studies!

Families with a higher number of cases than expected -> genetics are involved!

Family study - example from multiple sclerosis

	Risk Ratio (RR)
Monozygotic twin	23.62 (8.71-64.02)
Dizygotic twin	2.18 (0.71-6.68)
Sibling	7.13 (6.42-7.93)
Adopted sibling	1.87 (0.23-15.46)
Cousin	1.63 (1.36-1.97)

Westerlind et al, Brain 2014

1









GWAS - the how to!



On to the genetics...

2

Karolinska Institutet

Karolinska Institutet

Gut, 2017 Mar;66(3):421-428. doi: 10.1136/gut/inl-2015-309934. Epub 2015 Nov 2.

Dense genotyping of immune-related loci identifies HLA variants associated with increased risk of collagenous colitis.

Westerleich H¹², Malander MB¹⁴, Besson F^{23,4}, March A⁶, Rondoll E⁷, Assad G², Ratter J², Hilberhall M⁰, Lieb W⁷, Kälters H¹, Brinedel B¹, Patrialeou L³, Halfvanson J⁰, Totkvat L⁴, Block J⁴, Andreasson J⁰, Assaul L⁹, Amerika S¹⁴, Merrike S¹⁰, Madisch A¹¹, Ottisson B¹², Lötteric B^{3,13}, Hultzonic B¹⁴, Erinke A⁶, D'Arnalo M^{2,18}.

Collagenous colitis

- Inflammatory bowel disease
- Onset > 50 years
 Predominant in women
- i rodoninani ir won
- Chronic watery diarrhea
 Diagnosis through biopsies showing deposition of collagen in lamina propria
- Etiology unknown!
- · Reports about familial clustering, no proper genetic study done

Study material

- · 314 cases, 4,299 controls from Sweden and Germany
- · Genotyped on the Immunochip

Karolinska Institutet



AND INDIVIDUALS!



Genotyping arrays

Oligonucleotide probes

→single nucleotide extension by labeled nucleotides
→allele specific probe + labeled ss DNA fragments
→etc

Karolinska







Karolinska Institutet

QC of genotyping

Pipeline for QC SNP markers

v Table 1. OC marker: nineline //mmunochin.discouses)					
	SE SNPs dropped	SE SNPs	DE SNPs dropped	DE SNPs	
		196524		196524	
C (including intensity outliers) *	52762	143762	52762	143762	
	17174	126588	19495	124267	
8%	142	126446	78	124189	
Excelsion and excelsion of Sector 2014		405000	40.07	440070	

Success rate = 30%	142	120440	/0	124109
Differential missing cases vs controls	1118	125328	4937	119252
Minor allele frequency < 0.01	6737	118591	6078	113174
HWE P < 1.0E-07	87	118504	275	112899
Passing QC in both datasets	787	110634	2265	110634



QC step	SE cases *	SE ctris #	DE cases #	DE ctris
Initial number	163	2046	91	2018
Genotyping success rate < 95%	0	5	0	0
Genotype-phenotype sex mismatch	1	5	1	13
Genotype relatedness > 0.1	2	39	3	31
PCA outlier (> 6 SD)	0	2	0	0
Final sample size	160	1995	87	1974

- Done using the software PLINK and the command –indep-pairwise
- GWAS: 50 5 0.5
- ICHIP: 50 5 0.2*
- "Problematic" genomic regions with high LD and/or known inversion* must also be removed

*Abraham G, PLoS One, 2014



Odds ratio (recap I hope)

	Cases	Controls	Total
Exposed	а	c	a+c
Unexposed	b	d	b+d
Total	a+b	c+d	

$$OR = \frac{a/b}{c/d} = \frac{a*d}{c*d}$$







Human Leucocyte Antigen

HLA class I: Peptides from *inside* the cell to CD8+T-cells (killer T-cells)

HLA class II: Antigens from *outside* the cell to CD4+ T-cells (stimulating B-cells)



Most autoimmune disease have an HLA association

- Is collagenous collitis maybe an autoimmune disease?
- We want to go from SNP-data to HLA genotypes!



Imputation is the way to go!

A bi

Karolinska Institutet

A bit more about genetics first!

The central dogma



Outline

- The genome is not random
 → How the genome is structured
- Genetics differ between populations
 → A bit about migration and evolution of mankind
- Recombination creates diversity
 → More on inheritance
- Statistical modeling and software
 → The HowTo







A protein is not translated from a coherent stretch of DNA Linkage disequilibrium (LD) • It's beneficial to preserve function, so there is high(er) correlation within a (functional) stretch of DNA. • The LD differs across the genome. • If we know the LD structure, we can infere the genotype at position B for an individual if we know the contropt at position Al	
 It's beneficial to preserve function, so there is high(er) correlation within a (functional) stretch of DNA. The LD differs across the genome. If we know the LD structure, we can infere the genotype at position B for an individual if we know the enotype at position AI 	
Pic: news-medical.net	
Karolinska Institutet	nska
LD – more formal LD structure – an example	
• $D_{AB} = p_{AB} - p_A p_B$ • $D' = D'_{D_{min}}$ where $D_{min} = \left\{ \max(-p_{AB}, -(1-p_A)(1-p_B)) \text{ when } D < 0 \\ \min(p_A(1-p_B), (1-p_A)p_B) \text{ when } D > 0 \\ \cdot r = \frac{D}{\sqrt{p_A(1-p_A)p_B(1-p_B)}} \right\}$	
Pic: Olsan et al, BMC Genetics, 2007	_



The genetics differ between populations

Genetics differ across the world!

- Genetic variants differ in frequency between populations
- And certain variants might not even exist in some populations
- Addmixture, bottlenecks and selective pressure are some explanations for this

Karolinska Institutet



So somehow we need to take population into account ...

Mitosis vs meiosis

- Mitosis: the cell copies itself. These are diploid (46 chromosomes, 23 pairs)
- Meiosis: creation of gametes (eggs/sperms), which are haploid (23 chromosomes)



Recombination creates diversity

- Is more common at certain position in the genome, so called recombination hotspots
- Is part in driving evolution (together with mutations, genetic drift, system of mating, population structure, selection and genetic linkage)
- But can also create dysfunctional genotypes that might be fatal
- · By studying families we can learn more about recombination



Karolinska Institutet

In summary

- The correlation structure in the genome is called LD
- · The genetics differ between populations
- · There will be diversity in a population due to recombination

Karolinska Institutet



More on the statistics

We can estimate

 \rightarrow the correlation between different positions in the genome \rightarrow And the frequencies of the genotypes in the population

And we can use this to estimate the genotype at position i+1 conditional on position i



The genome is a Markov model!

(which is outside the scope of this course, but you should at least know there's a statistical way to model the genome)





- A key publication is Lie & Stephens, Genetics 2003 (HOTSPOTTER)
- · Lie and Stephens' approach is used by several software (BEAGLE, IMPUTE, MaCH)
- By
 → assuming a Markov model
 - → estimating the LD structure
 → and the genotypic frequencies
 - we can infere the missing genotypes
- A key point is (of course) to have an accurate reference panel.
- Remember that mutations and recombinations introduce uncertatinty to the imputation.
- · But most software also give a quality statistic of the imputed genotypes.



Karolinska Institutet

So by using imputation, we can go from SNP-data to genotypes

(there are other ways of doing this, this is just one version of inferring genotypes, but it's the gold standard for estimating the genotypes in the HLA region).

	TT	CT	CC	RAF	p Value	OR (95% CI)
Discovery						
Cases	8	94	145	0.222	1.6×10 ⁻⁸	1.98 (1.56 to 2.51)
Controls	59	895	3015	0.128		
Replication						
Cases	3	24	40	0.224	1.3×10 ⁻⁴	2.79 (1.65 to 4.73
Controls	4	63	263	0.108		
Combined						
Cases	11	118	185	0.223	2.3×10*11	2.06 (1.67 to 2.55
Controls	63	958	3278	0.126		

Westerlind et al, GUT 2017



In conclusion

Collagenous colitis seem to have an association to HLA class II, implicating that it's an autoimmune disease!



Summary of the lecture

- Quality control of genetic data is important!
- · We need to know about genetics and evolution to be able to use the correct statistics
- Statistics is just statistics, and inferring causality is not trivial. Association does not imply causation!!
- · We need a plausible explanation, replication and preferably also functional studies to trust the results fully!