

B I O I N F O R M A T I C S

Kristel Van Steen, PhD²

Montefiore Institute - Systems and Modeling

GIGA - Bioinformatics

ULg

kristel.vansteen@ulg.ac.be

CHAPTER 6: POPULATION BASED ASSOCIATION STUDIES

1 Introduction

1.a Dissecting human disease in the postgenomic era

1.b Genetic association studies

2 Preliminary analyses

2.a Hardy-Weinberg equilibrium

2.b Missing genotype data

2.c Haplotype and genotype data

2.d Measures of LD and estimates of recombination rates

2.e SNP tagging

3 Tests of association: single SNP

4 Tests of association: multiple SNPs

5 Dealing with population stratification

5.a Spurious associations

5.b Genomic control

5.c Structured association methods

5.d Other approaches

6 Multiple testing

6.a General setting

6.b Controlling the type I error

7 Assessing the function of genetic variants

8 Proof of concept

1 Introduction

1.a Dissecting human disease in the postgenomic era

THE HUMAN GENOME

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹ Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Yandell,¹ Charyl A. Brans,¹ Robert A. Holt,¹ Jeanine D. Gocayne,¹ Peter Anagnostides,¹ Richard M. Ballow,¹ David H. Hays,¹ Jennifer Russo-Wertman,¹ Qing Zhang,¹ Chinnappa D. Kodira,¹ Xiangqun H. Zhang,¹ Liu Chan,¹ Martin Skupski,¹ Chagadharan Subramanian,¹ Paul D. Thomas,¹ Jinghui Zhang,¹ George L. Gabor Miklos,¹ Catharina Nelson,¹ Susmita Broder,¹ Andrew C. Clark,¹ Joe Hradeau,¹ Victor A. Holtkamp,¹ Marton Zeller,¹ Arnold J. Levine,¹ Richard J. Roberts,¹ Paul Simon,¹ Carolyn Stanyon,¹ Michael Hunkapiller,¹ Randall Bolanos,¹ Arthur Delcher,¹ Ben Dawe,¹ Daniel Fasilo,¹ Michael Flanagan,¹ Ulises Flores,¹ Aaron Hahn,¹ Sriharu Hannanholi,¹ Saul Kravitz,¹ Samuel Levy,¹ Clark Mobarry,¹ Kaut Balalet,¹ Karin Ramington,¹ Jane Abu-Thalab,¹ Ellen Baskey,¹ Kendra Biddick,¹ Vivian Bonazzi,¹ Rhonda Brandon,¹ Michele Caselli,¹ Ishwar Chandrasekharan,¹ Rosana Charlab,¹ Kishi Chikara,¹ Zhenping Dong,¹ Valerio Di Pasquale,¹ Patrick Dowd,¹ Steven Ellsack,¹ Carlos Brangelato,¹ Andrei E. Cobralian,¹ Waihu Cao,¹ Wangmao Guo,¹ Fengcheng Gong,¹ Zhenping Guo,¹ Ping Guan,¹ Thomas J. Hainan,¹ Maureen E. Higgins,¹ Rui-Ru Ji,¹ Zhixi Ke,¹ Karen A. Ketchum,¹ Zhonghua Li,¹ Yiding Liu,¹ Zherya Li,¹ Jiyin Li,¹ Yong Liang,¹ Xiaoying Lin,¹ Fu Lu,¹ Gennady V. Markovits,¹ Natalia Pichshina,¹ Hulan H. Moore,¹ Ashwinkumar K. Nair,¹ Václav A. Hozayin,¹ Beate Muehlen,¹ Deborah Husskara,¹ Douglas S. Rasch,¹ Steven Salzberg,¹ Wai Shao,¹ Bkiong Shue,¹ Jingtao Sun,¹ Zhen Yuan Wang,¹ Aihui Wang,¹ Xin Wang,¹ Jian Wang,¹ Ming-Mei Wu,¹ Ron Wides,¹ Chunlin Xiao,¹ Chunhua Yan,¹ Alison Yao,¹ Jane Yu,¹ Ming Zhan,¹ Weiqiang Zhang,¹ Hongyi Zhang,¹ Qi Zhao,¹ Liandeng Zheng,¹ Fei Zhong,¹ Weyan Zhong,¹ Shiqing C. Zhu,¹ Shuying Zhuo,¹ Dennis Gilbert,¹ Suzanna Bammertner,¹ Gene Spier,¹ Christine Carter,¹ Anibal Cavalli,¹ Trevor Woodgates,¹ Farooq Ali,¹ Hejira An,¹ Adarshika An,¹ Danika Baldwin,¹ Holly Baden,¹ Mary Barstead,¹ Ian Barrow,¹ Karen Beeson,¹ Diana Busam,¹ Amy Carver,¹ Angela Center,¹ Ming Lai Cheng,¹ Liz Curry,¹ Steve Danaher,¹ Lionel Davenport,¹ Raymond Deshats,¹ Susanne Dietz,¹ Kristina Dodson,¹ Lisa Dool,¹ Steven Ferreira,¹ Neha Garg,¹ Andrea Giesemann,¹ Brett Hart,¹ Jason Higgins,¹ Charles Hynes,¹ Cheryl Heiner,¹ Suzanne Hudson,¹ Doumon Hostin,¹ Jarrett Hoock,¹ Timothy Howland,¹ Chinyere Ikegwam,¹ Jeffery Johnson,¹ Francis Kalush,¹ Lesley Klein,¹ Shashi Koduru,¹ Amy Love,¹ Felicia Mann,¹ David May,¹ Steven McCawley,¹ Tina McIntosh,¹ Ivy McMillan,¹ Moe Moy,¹ Linda Moy,¹ Brian Murphy,¹ Keith Nelson,¹ Cynthia Frankoch,¹ Eric Pratt,¹ Vikita Puri,¹ Hans Qawash,¹ Matthew Raardon,¹ Robert Rodriguez,¹ Yasuki Rogers,¹ Osamu Rosenthal,¹ Bob Rubin,¹ Richard Scott,¹ Cynthia Sider,¹ Michelle Smallwood,¹ Eric Stewart,¹ Renee Strong,¹ Ellen Suh,¹ Reginald Thomas,¹ Ni Ni Tint,¹ Sukyee Tse,¹ Claire Vech,¹ Gary Wang,¹ Jeremy Wetter,¹ Sharita Williams,¹ Monica Williams,¹ Sandra Windsor,¹ Emily Winn-Dean,¹ Karla Wolk,¹ Jaydree Zaveri,¹ Karen Zaveri,¹ Josep F. Abril,¹ Roderic Aiguo,¹ Michael J. Campbell,¹ Kimman V. Sjostrand,¹ Brian Karbak,¹ Anshu Kojima,¹ Tatsuya Mi,¹ Barry Lazarus,¹ Thomas Hutton,¹ Aparna Hanchahal,¹ Karen Olaner,¹ Anahya Muruganujan,¹ Nan Guo,¹ Shiji Sato,¹ Vinat Batra,¹ Sorin Istrail,¹ Ross Lippert,¹ Russell Schwartz,¹ Brian Walenz,¹ Shibu Yooseph,¹ David Allan,¹ Anand Basal,¹ James Bonandale,¹ Louis Blöck,¹ Marcelo Caminha,¹ John Carnas-Silva,¹ Parviz Caslik,¹ Yue-Hui Chiang,¹ My Coyma,¹ Carl Dahlke,¹ Anja Dashtartayeva,¹ Maria Dombocki,¹ Michael Donnelly,¹ Dalia El,¹ Shiva Eppanham,¹ Carl Foster,¹ Harold Gilo,¹ Stephen Glanville,¹ Kamath Gleser,¹ Anna Glöckle,¹ Mark Gorkov,¹ Kee Graham,¹ Barry Gropman,¹ Michael Harris,¹ Jeremy Hall,¹ Scott Henderson,¹ Jeffrey Hoover,¹ Donald Jennings,¹ Catharine Jordan,¹ James Jordan,¹ John Kachra,¹ Leonid Kagan,¹ Cheryl Kraft,¹ Alexander Lewitsky,¹ Mark Lewis,¹ Xianglin Liu,¹ John Lopez,¹ Daniel Ma,¹ William Majoros,¹ Joe McGinnis,¹ Sean Murphy,¹ Matthew Newman,¹ Young Nguyen,¹ Hyeon Nguyen,¹ Marc Odell,¹ Sun Park,¹ Jim Park,¹ Maratli Petersen,¹ William Rhee,¹ Robert Sanders,¹ John Scott,¹ Michael Simpson,¹ Thomas Smith,¹ Arisa Sprague,¹ Timothy Stockwell,¹ Russell Turner,¹ Eli Venter,¹ Hai Wang,¹ Melyuan Wan,¹ David Wu,¹ Mitchell Wu,¹ Ashley Xia,¹ Ali Zandieh,¹ Xiaohong Zhu,¹

16 FEBRUARY 2001 VOL 29 | SCIENCE www.science.org

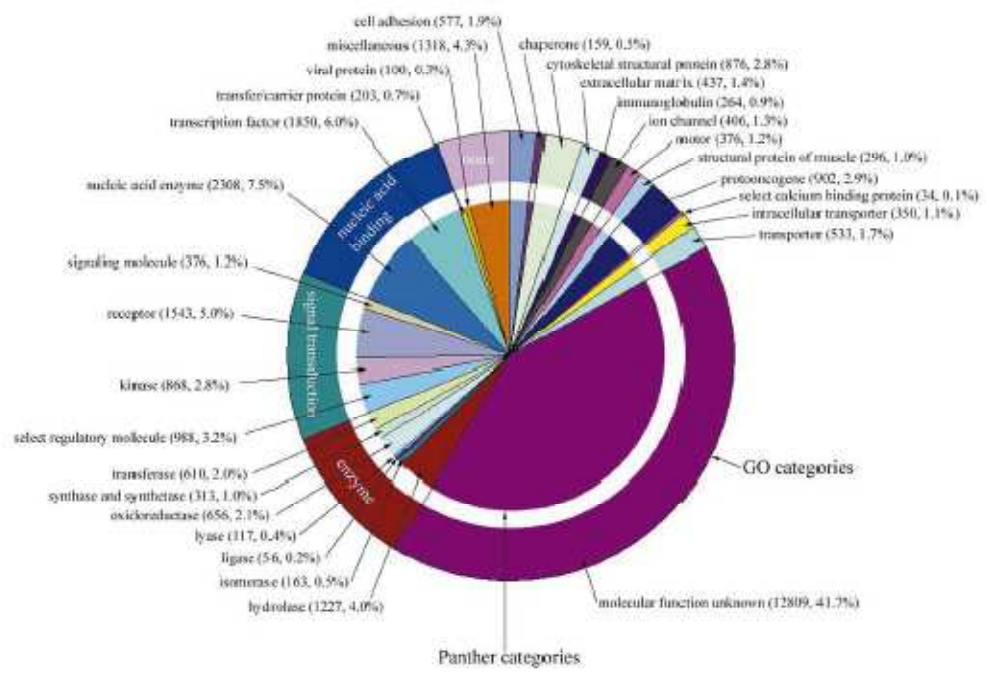


Fig. 15. Distribution of the molecular functions of 26,383 human genes. Each slice lists the numbers and percentages (in parentheses) of human gene functions assigned to a given category of molecular function. The outer circle shows the assignment to molecular function categories in the Gene Ontology (GO) (179), and the inner circle shows the assignment to Celera's Panther molecular function categories (176).

Introduction

- The complete genome sequence of humans and of many other species provides a new starting point for understanding our basic genetic makeup and how variations in our genetic instructions result in disease.
- The pace of the molecular dissection of human disease can be measured by looking at the catalog of human genes and genetic disorders identified so far in *Mendelian Inheritance in Man* and in *OMIM*, its online version, which is updated daily (www.ncbi.nlm.nih.gov/omim).

(V. A. McKusick, *Mendelian Inheritance in Man* (Johns Hopkins Univ. Press, Baltimore, ed. 12, 1998))

Introduction

NCBI

OMIM
Online Mendelian Inheritance in Man

Johns Hopkins University

My NCBI [Sign In] [Register]

All Databases PubMed Nucleotide Protein Genome Structure PMC OMIM

Search OMIM for

Entrez

OMIM

Search OMIM
Search Gene Map
Search Morbid Map

Help
OMIM Help
How to Link

FAQ
Numbering System
Symbols
How to Print
Citing OMIM
Download

OMIM Facts
Statistics

Limits Preview/Index History Clipboard Details

- Enter one or more search terms.
- Use **Limits** to restrict your search by search field, chromosome, and other criteria.
- Use **Index** to browse terms found in OMIM records.
- Use **History** to retrieve records from previous searches, or to combine searches.

OMIM® - Online Mendelian Inheritance in Man®

Welcome to OMIM®, Online Mendelian Inheritance in Man®. OMIM is a comprehensive, authoritative, and timely compendium of human genes and genetic phenotypes. The full-text, referenced overviews in OMIM contain information on all known mendelian disorders and over 12,000 genes. OMIM focuses on the relationship between phenotype and genotype. It is updated daily, and the entries contain copious links to other genetics resources.

This database was initiated in the early 1960s by Dr. Victor A. McKusick as a catalog of mendelian traits and disorders, entitled Mendelian Inheritance in Man (MIM). Twelve book editions of MIM were published between 1966 and 1998. The online version, OMIM, was created in 1985 by a collaboration between the National Library of Medicine and the William H. Welch Medical Library at Johns Hopkins. It was made generally available on the internet starting in 1987. In 1995, OMIM was developed for the World Wide Web by NCBI, the National Center for Biotechnology Information.

OMIM is authored and edited at the McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine.

Introduction

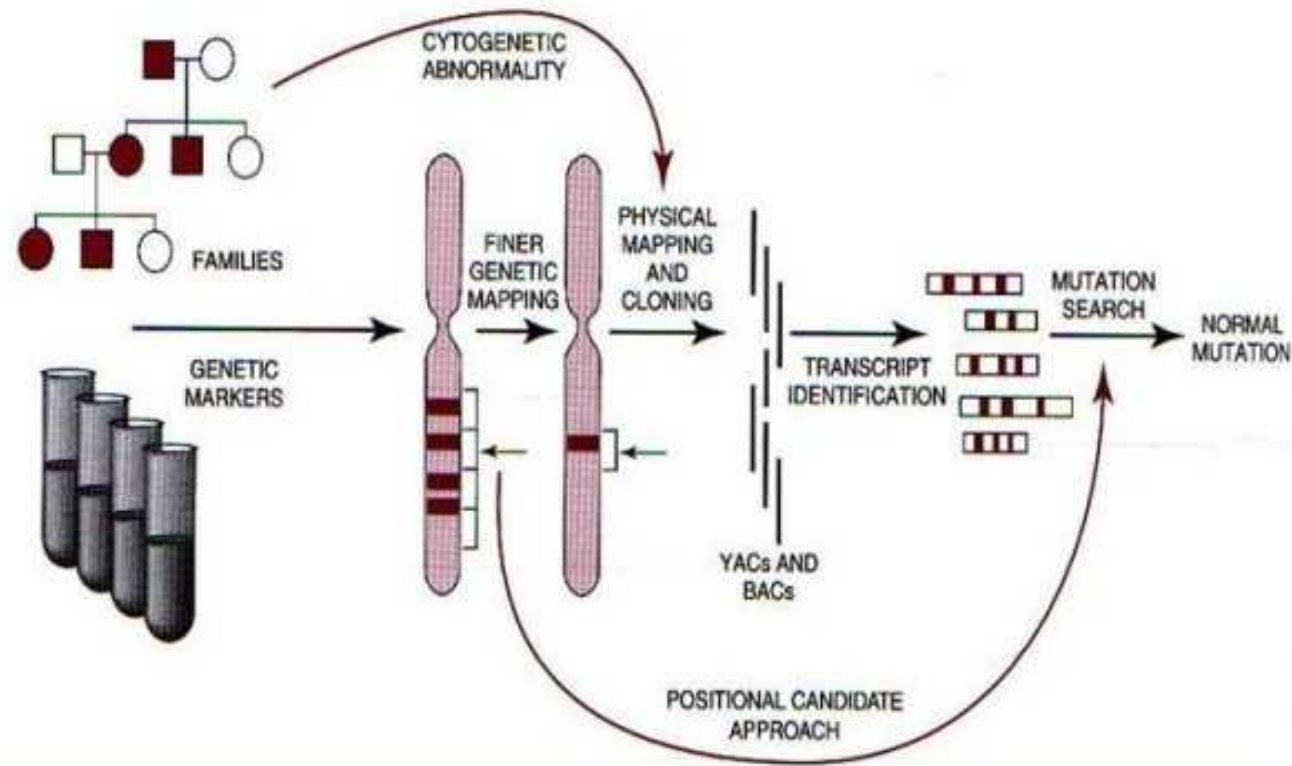
- OMIM Statistics for November 1, 2009: Number of entries

	Autosomal	X-Linked	Y-Linked	Mitochondrial	Total
* Gene with known sequence	<u>12266</u>	<u>604</u>	<u>48</u>	<u>35</u>	<u>12953</u>
+ Gene with known sequence and phenotype	<u>331</u>	<u>21</u>	0	<u>2</u>	<u>354</u>
# Phenotype description, molecular basis known	<u>2400</u>	<u>214</u>	<u>4</u>	<u>26</u>	<u>2644</u>
≈ Mendelian phenotype or locus, molecular basis unknown	<u>1646</u>	<u>141</u>	<u>5</u>	0	<u>1792</u>
Other, mainly phenotypes with suspected mendelian basis	<u>1874</u>	<u>137</u>	<u>2</u>	0	<u>2013</u>
Total	<u>18517</u>	<u>1117</u>	<u>59</u>	<u>63</u>	<u>19756</u>

Introduction

- Beginning in 1986, map-based gene discovery (positional cloning) became the leading method for elucidating the molecular basis of genetic disease.
- Almost all medical specialties have used this approach to identify the genetic causes of some of the most puzzling human disorders.
- Positional cloning has also been used reasonably successfully in the study of common diseases with multiple causes (so-called *complex disorders*), such as type I diabetes mellitus and asthma.

Positional cloning



(http://www.molecularlab.it/public/data/GFPina/200924223125_positional%20cloning.JPG)

Terminology

- **BAC:** Bacterial Artificial Chromosome. A type of cloning vector derived from the naturally-occurring F factor episome. A BAC can carry 100 - 200 kb of foreign DNA / YAC: Yeast Artificial Chromosome
- **Cloning vector:** A DNA construct capable of replication within a bacterial or yeast host that can harbor foreign DNA, facilitating experimental manipulation of that DNA segment.
- **Complex disease:** Condition caused by many contributing factors. Such a disease is also called a multifactorial disease.
 - Some disorders, such as sickle cell anemia and cystic fibrosis, are caused by mutations in a single gene.
 - Common medical problems such as heart disease, diabetes, and obesity likely associated with the effects of multiple genes in combination with lifestyle and environmental factors.

Introduction

- With the availability of the human genome sequence and those of an increasing number of other species, sequence-based gene discovery is complementing and will eventually replace map-based gene discovery.
- These and other recent developments in the field have caused a paradigm shift in biomedical research:

Structural genomics	→	Functional genomics
Genomics	→	Proteomics
Map-based gene discovery	→	Sequence-based gene discovery
Monogenic disorders	→	Multifactorial disorders
Specific DNA diagnosis	→	Monitoring of susceptibility
Analysis of one gene	→	Analysis of multiple genes in gene families, pathways, or systems
Gene action	→	Gene regulation
Etiology (specific mutation)	→	Pathogenesis (mechanism)
One species	→	Several species

Introduction

- Initial analyses of the completed chromosomal sequences suggest that the number of human genes is lower than expected.
- These findings are consistent with the idea that variations in gene regulation and the splicing of gene transcripts explain how one protein can have distinct functions in different types of tissue.
- At the beginning of the 21st century, it also seemed likely that obvious deleterious mutations in the coding sequences of genes are responsible for only a fraction of the differences in disease susceptibility between individuals, and that sequence variants affecting gene splicing and regulation must play an important part in determining disease susceptibility.

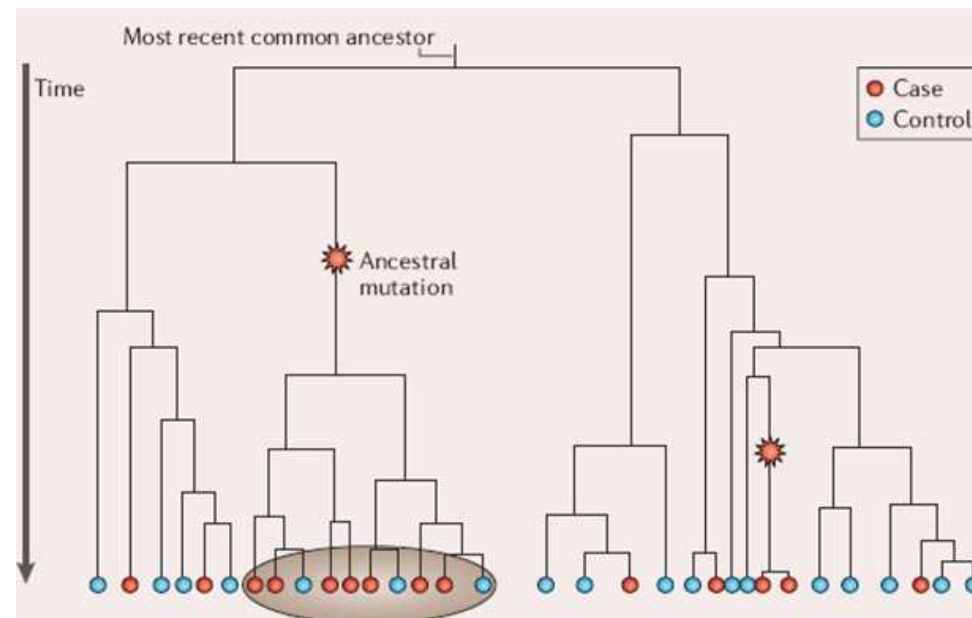
Introduction

- As only a small proportion of the millions of sequence variations in our genomes will have such functional impacts, identifying this subset of sequence variants is a challenging task.
- The success of global efforts to identify and annotate sequence variations in the human genome, which are called single-nucleotide polymorphisms (SNPs), is reflected in the abundance of SNP databases
 - www.ncbi.nlm.nih.gov/SNP,
 - <http://snp.cshl.org>,
 - <http://hgbase.cgr.ki.se>.
- However, the follow-up work of understanding how these and other genetic variations regulate the phenotypes (visual characteristics) of human cells, tissues, and organs will occupy biomedical researchers for all of the 21st century

1.b Population-based genetic association studies

Introduction

- The goal of population association studies is to identify patterns of polymorphisms that vary systematically between individuals with different disease states and could therefore represent the effects of risk-enhancing or protective alleles.



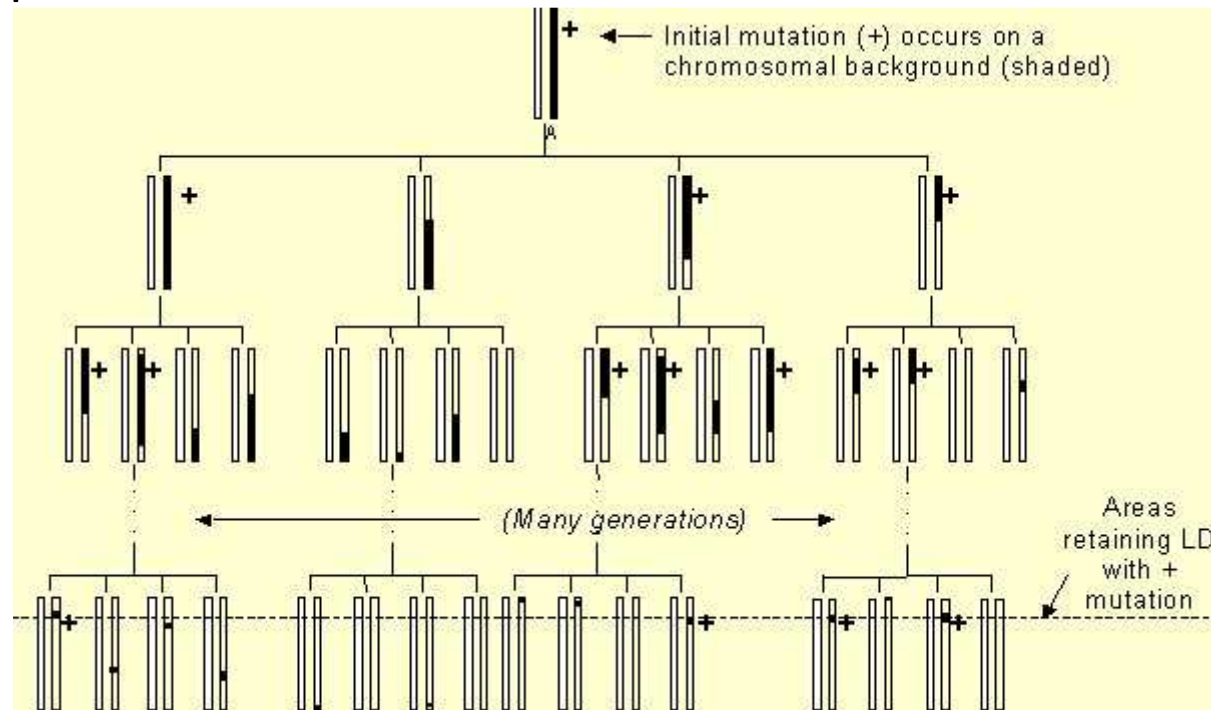
(Balding 2006)

Introduction

- When performing a genetic association study, there are a number of pitfalls one should be aware of.
- Perhaps the most crucial one is related to the realization that some patterns may arise simply by chance.
- To distinguish between true and chance effects, there are two routes to be taken:
 - Set tight standards for statistical significance
 - Only consider patterns of polymorphisms that could plausibly have been generated by causal genetic variants (use understanding of human genetic history or evolutionary processes such as recombination or mutation)
 - Adequately deal with distorting factors, including missing data and genotyping errors (quality control measures)

Introduction

- Hence, the key concept in a (population-based) genetic association study is linkage disequilibrium.



- This gives the rationale for performing genetic association studies

Types of genetic association studies

- Candidate polymorphism
 - These studies focus on an individual polymorphism that is suspected of being implicated in disease causation.
- Candidate gene
 - These studies might involve typing 5–50 SNPs within a gene (defined to include coding sequence and flanking regions, and perhaps including splice or regulatory sites).
 - The gene can be either a positional candidate that results from a prior linkage study, or a functional candidate that is based, for example, on homology with a gene of known function in a model species.

Types of genetic association studies

- Fine mapping
 - Often refers to studies that are conducted in a candidate region of perhaps 1–10 Mb and might involve several hundred SNPs.
 - The candidate region might have been identified by a linkage study and contain perhaps 5–50 genes.
- Genome-wide
 - These seek to identify common causal variants throughout the genome, and require $\geq 300,000$ well-chosen SNPs (more are typically needed in African populations because of greater genetic diversity).
 - The typing of this many markers has become possible because of the International HapMap Project and advances in high-throughput genotyping technology

Types of population association studies

- The aforementioned classifications are not precise: some candidate-gene studies involve many hundreds of genes and are similar to genome-wide scans.
- Typically, a causal variant will not be typed in the study, possibly because it is not a SNP (it might be an insertion or deletion, inversion, or copy-number polymorphism).
- Nevertheless, a well-designed study will have a good chance of including one or more SNPs that are in strong linkage disequilibrium with a common causal variant.

Analysis of population association studies

- Statistical methods that are used in pharmacogenetics are similar to those for disease studies, but the phenotype of interest is drug response (efficacy and/or adverse side effects).
- In addition, pharmacogenetic studies might be prospective whereas disease studies are typically retrospective.
- Prospective studies are generally preferred by epidemiologists, and despite their high cost and long duration some large, prospective cohort studies are currently underway for rare diseases.
- Often a case–control analysis of genotype data is embedded within these studies, so many of the statistical analyses that are discussed in this chapter can apply both to retrospective and prospective studies.
- However, specialized statistical methods for time-to-event data might be required to analyse prospective studies.

Analysis of population association studies

- Design issues guide the analysis methods to choose from:

	Details	Advantages	Disadvantages	Statistical analysis method
Cross-sectional	Genotype and phenotype (ie, note disease status or quantitative trait value) a random sample from population	Inexpensive. Provides estimate of disease prevalence	Few affected individuals if disease rare	Logistic regression, χ^2 tests of association or linear regression
Cohort	Genotype subsection of population and follow disease incidence for specified time period	Provides estimate of disease incidence	Expensive to follow-up. Issues with drop-out	Survival analysis methods
Case-control	Genotype specified number of affected (case) and unaffected (control) individuals. Cases usually obtained from family practitioners or disease registries, controls obtained from random population sample or convenience sample	No need for follow-up. Provides estimates of exposure effects	Requires careful selection of controls. Potential for confounding (eg, population stratification)	Logistic regression, χ^2 tests of association
Extreme values	Genotype individuals with extreme (high or low) values of a quantitative trait, as established from initial cross-sectional or cohort sample	Genotype only most informative individuals hence save on genotyping costs	No estimate of true genetic effect sizes	Linear regression, non-parametric, or permutation approaches
Case-parent triads	Genotype affected individuals plus their parents (affected individuals determined from initial cross-sectional, cohort, or disease-outcome based sample)	Robust to population stratification. Can estimate maternal and imprinting effects	Less powerful than case-control design	Transmission/disequilibrium test, conditional logistic regression or log-linear models
Case-parent-grandparent septets	Genotype affected individuals plus their parents and grandparents	Robust to population stratification. Can estimate maternal and imprinting effects	Grandparents rarely available	Log-linear models
General pedigrees	Genotype random sample or disease-outcome based sample of families from general population. Phenotype for disease trait or quantitative trait	Higher power with large families. Sample may already exist from linkage studies	Expensive to genotype. Many missing individuals	Pedigree disequilibrium test, family-based association test, quantitative transmission/disequilibrium test
Case-only	Genotype only affected individuals, obtained from initial cross-sectional, cohort, or disease-outcome based sample	Most powerful design for detection of interaction effects	Can only estimate interaction effects. Very sensitive to population stratification	Logistic regression, χ^2 tests of association
DNA-pooling	Applies to variety of above designs, but genotyping is of pools of anywhere between two and 100 individuals, rather than on an individual basis	Potentially inexpensive compared with individual genotyping (but technology still under development)	Hard to estimate different experimental sources of variance	Estimation of components of variance

Table 2: Study designs for genetic association studies

(Cordell and Clayton, 2005)

Analysis of population association studies

- The design of a genetic association study may refer to
 - subject design (see before)
 - marker design:
 - Which markers are most informative? Microsatellites? SNPs? CNVs?
 - Which platform is the most promising?
 - study scale:
 - Genome-wide
 - Genomic

Analysis of population association studies

- Marker design

- Recombinations that have occurred since the most recent common ancestor of the group at the locus can break down associations of phenotype with all but the most tightly linked marker alleles.
- This permits fine mapping if marker density is sufficiently high (say, ≥ 1 marker per 10 kb).
- When the mutation entered into the population a long time ago, then a lot of recombination processes may have occurred, and hence the haplotype harboring the disease mutation may be very small.
 - This favors typing a lot of markers and generating dense maps
 - The drawback is the computational and statistical burden involved with analyzing such huge data sets.

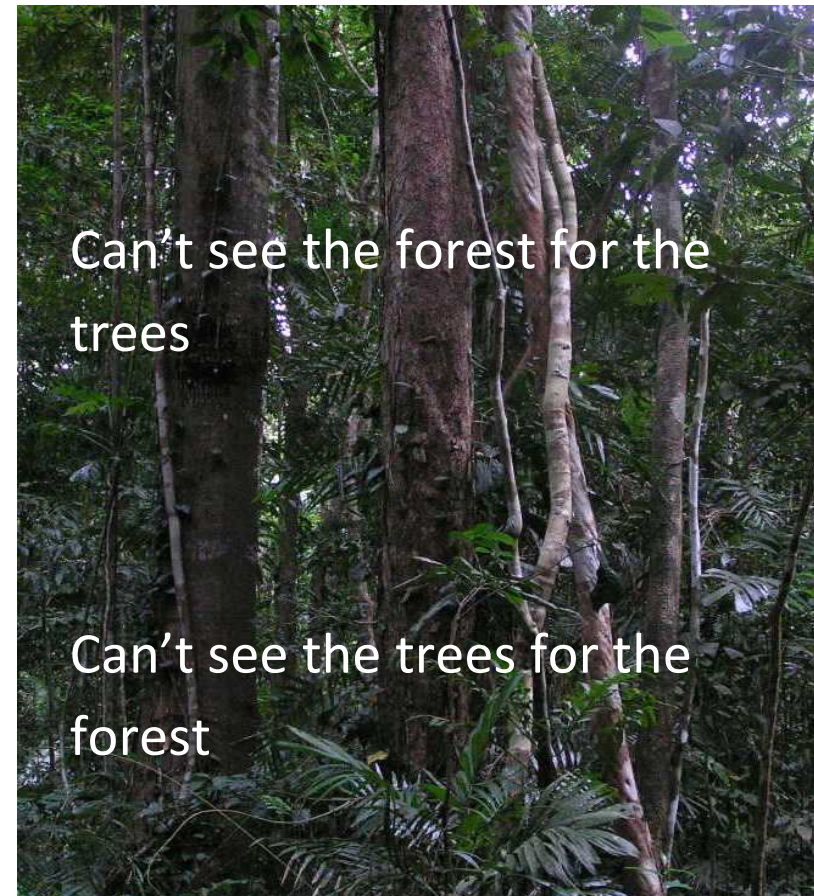
Analysis of population association studies

- Scale of genetic association studies

candidate gene approach

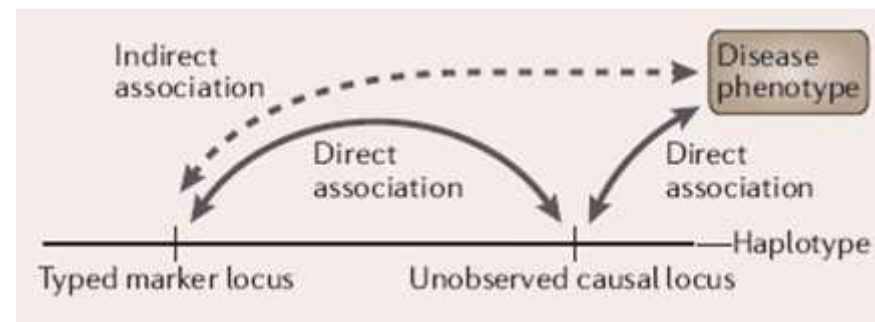
vs

genome-wide screening approach



Analysis of population association studies

- Direct versus indirect associations
 - The two direct associations that are indicated in the figure below, between a typed marker locus and the unobserved causal locus, cannot be observed, but if r^2 (a measure of allelic association) between the two loci is high then we might be able to detect the indirect association between marker locus and disease phenotype.



Power of genetic association studies

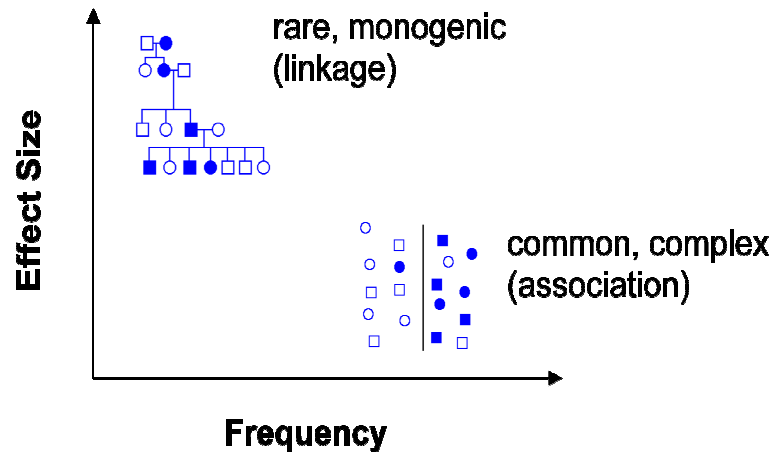
- Broadly speaking, association studies are sufficiently powerful only for common causal variants.
- The threshold for *common* depends on sample and effect sizes as well as marker frequencies, but as a rough guide the minor-allele frequency might need to be above 5%.
- The *common disease / common variant* (CDCV) hypothesis argues that genetic variations with appreciable frequency in the population at large, but relatively low ‘penetrance’ (or the probability that a carrier of the relevant variants will express the disease), are the major contributors to genetic susceptibility to common diseases.
 - If multiple rare genetic variants were the primary cause of common complex disease, association studies would have little power to detect them; particularly if allelic heterogeneity existed.
 - The major proponents of the CDCV were the movers and shakers behind the HapMap and large-scale association studies

Power of genetic association studies

- The competing hypothesis is cleverly the *Common Disease-Rare Variant* (CDRV) hypothesis. It argues that multiple rare DNA sequence variations, each with relatively high penetrance, are the major contributors to genetic susceptibility to common diseases.
- Although some common variants that underlie complex diseases have been identified, and given the recent huge financial and scientific investment in GWA, there is no longer a great deal of evidence in support of the CDCV hypothesis and much of it is equivocal...
- Both CDCV and CDRV hypotheses have their place in current research efforts.

Power of genetic association studies

Which gene hunting method is most likely to give success?



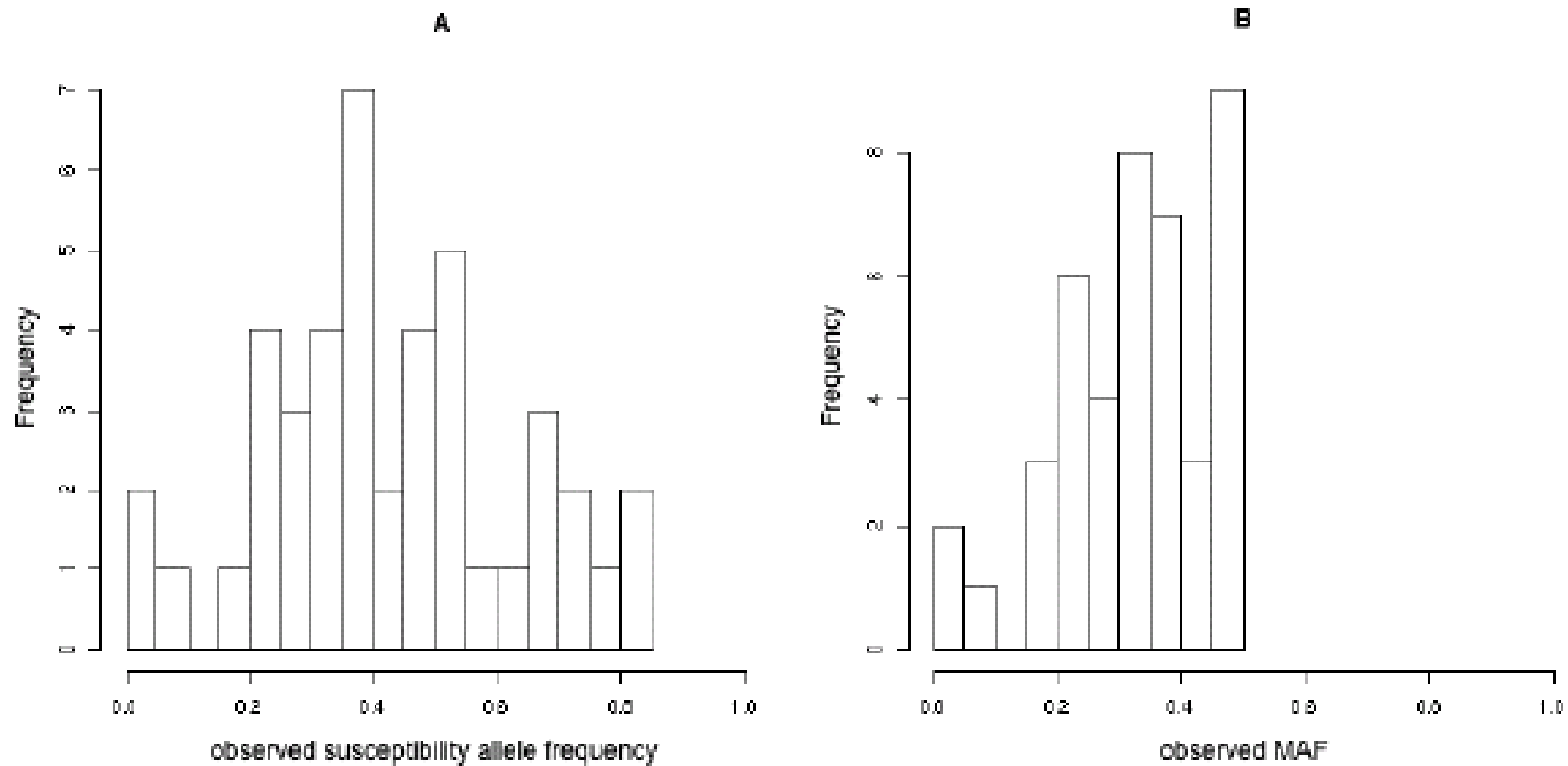
- Monogenic “Mendelian” diseases
 - Rare disease
 - Rare variants
 - Highly penetrant
- Complex diseases
 - Rare/common disease
 - Rare/common variants
 - Variable penetrance

(Slide: courtesy of Matt McQueen)

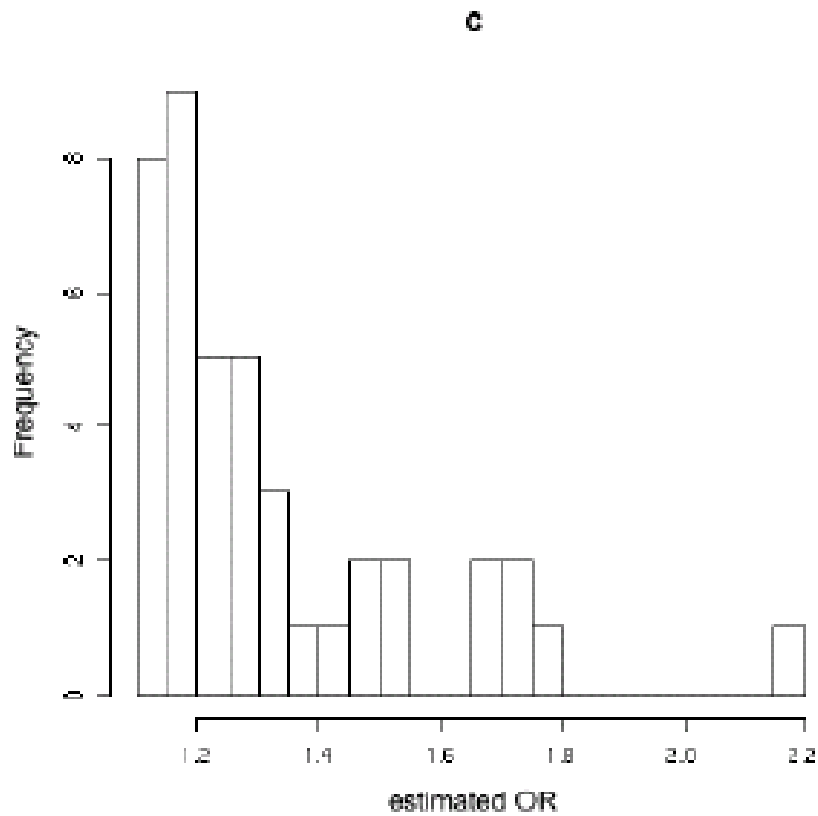
Power of genetic association studies

- Many genome scientists are turning back to study rare disorders that are traceable to defects in single genes, and whose causes have remained a mystery.
- The change is partly a result of frustration with the disappointing results of genome-wide association studies (GWAS). Rather than sequencing whole genomes, GWAS studies examine a subset of DNA variants in thousands of unrelated people with common diseases. Now, however, sequencing costs are dropping, and whole genome sequences can quickly provide in-depth information about individuals, enabling scientists to locate genetic mutations that underlie rare diseases by sequencing a handful of people.

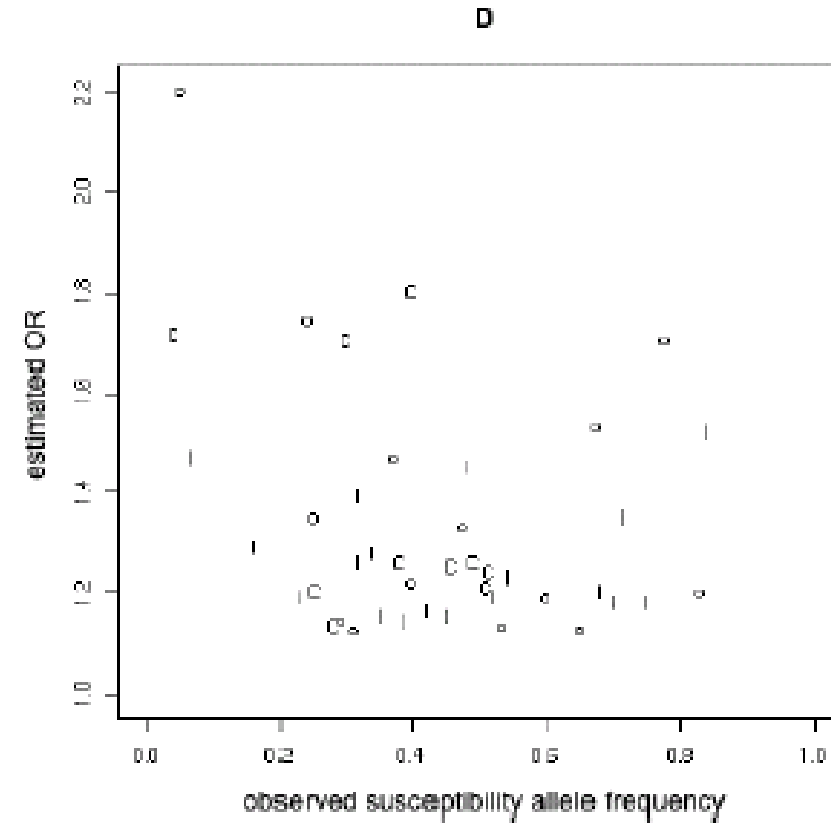
(Nature News: Published online 22 September 2009 | **461**, 459 (2009) | doi:10.1038/461458a)



(A and B) Histograms of susceptibility allele frequency and MAF, respectively, at confirmed susceptibility loci.



doi:10.1371/journal.pgen.0040033.g001



(C) Histogram of estimated ORs (estimate of genetic effect size) at confirmed susceptibility loci. (D) Plot of estimated OR against susceptibility allele frequency at confirmed susceptibility loci. (Iles 2008)

Factors influencing consistency of gene-disease associations

- Variables affecting inferences from experimental studies:
 - In vitro or in vivo system studied
 - Cell type studied
 - Cultured versus fresh cells studied
 - Genetic background of the system
 - DNA constructs
 - DNA segments that are included in functional (for example, expression) constructs
 - Use of additional promoter or enhancer elements
 - Exposures
 - Use of compounds that induce or repress expression
 - Influence of diet or other exposures on animal studies

(Rebbeck et al 2004)

Factors influencing consistency of gene-disease associations

- Variables affecting epidemiological inferences:
 - Inclusion/exclusion criteria for study subject selection
 - Sample size and statistical power
 - Candidate gene choice
 - A biologically plausible candidate gene
 - Functional relevance of the candidate genetic variant
 - Frequency of allelic variant
 - Statistical analysis
 - Consideration of confounding variables, including ethnicity, gender or age.
 - Whether an appropriate statistical model was applied (for example, were interactions considered in addition to main effects of genes?)
 - Violation of model assumptions

(Rebbeck et al 2004)

2 Preliminary analyses

2.a Introduction

2.b Hardy-Weinberg equilibrium

2.c Missing genotype data

2.d Haplotype and genotype data

Measures of LD and estimates of recombination rates

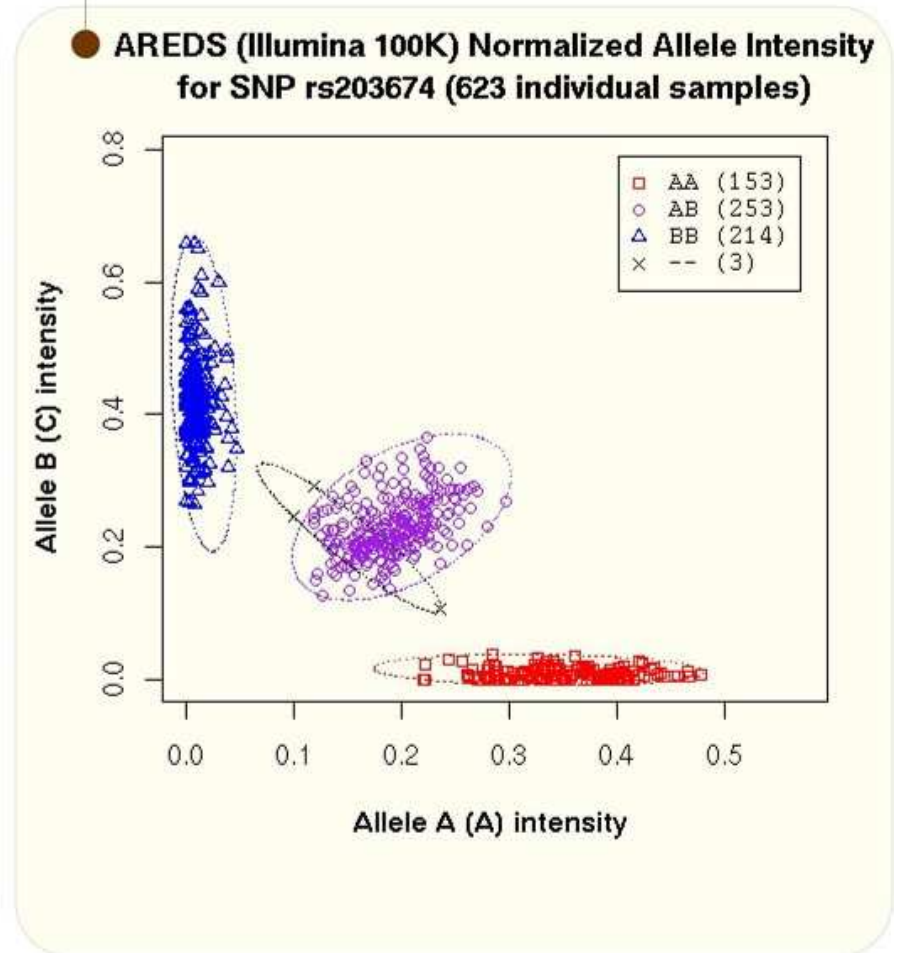
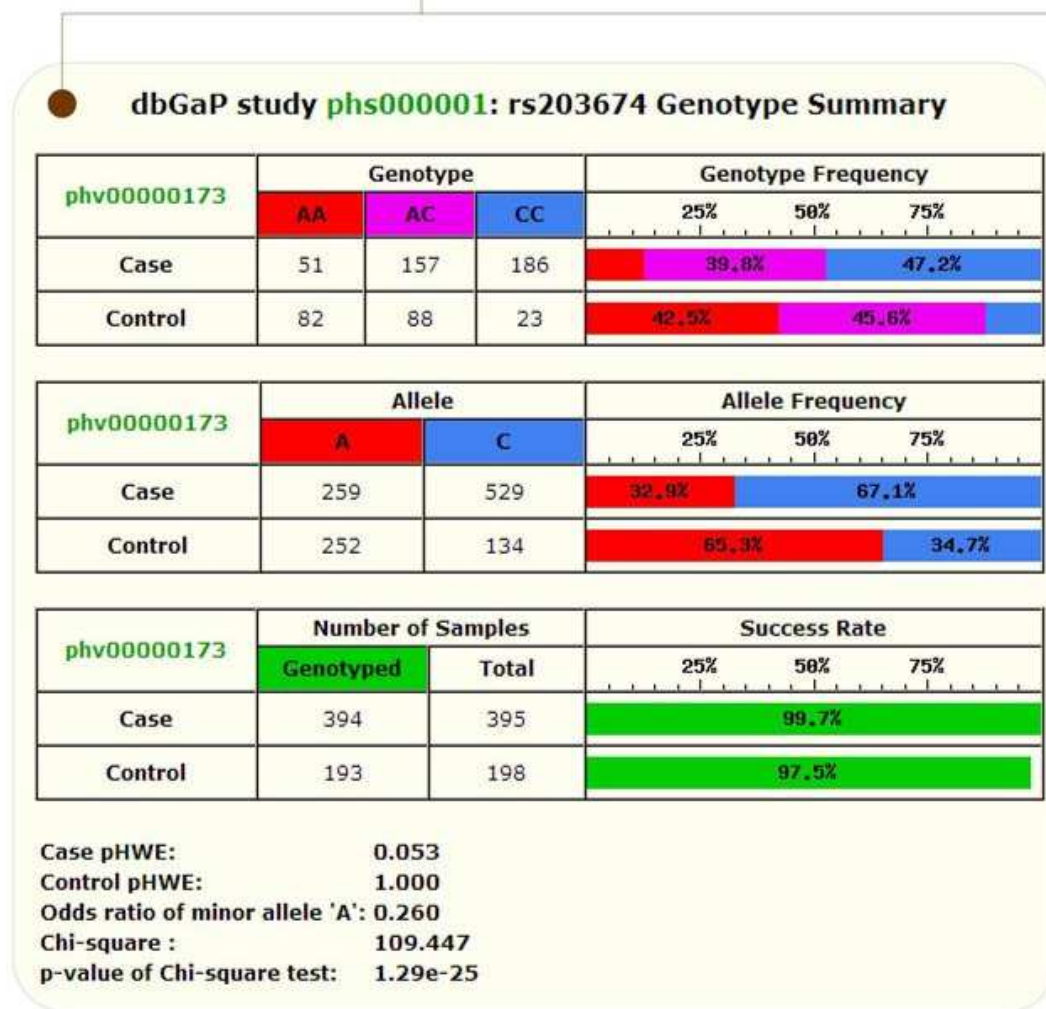
2.e SNP tagging

2.a Introduction

- Pre-analysis techniques often performed include:
 - testing for Hardy–Weinberg equilibrium (HWE)
 - strategies to select a good subset of the available SNPs ('tag' SNPs)
 - inferring haplotypes from genotypes.
- Data quality is of paramount importance, and data should be checked thoroughly before other analyses are started.
- Data should be checked for
 - batch or study-centre effects,
 - for unusual patterns of missing data,
 - for genotyping errors.

Introduction

- Recall that genotype data are not raw data:
 - Genotypes have been derived from raw data using particular software tools, one being more sensitive than the other
- For instance, SNP quality control involves assessing
 - missing data rates,
 - Hardy-Weinberg equilibrium (HWE),
 - allele frequencies,
 - Mendelian inconsistencies (using family-data)
 - sample heterozygosity, ...



(using dbGaP association browser tools)

2.b Hardy-Weinberg equilibrium

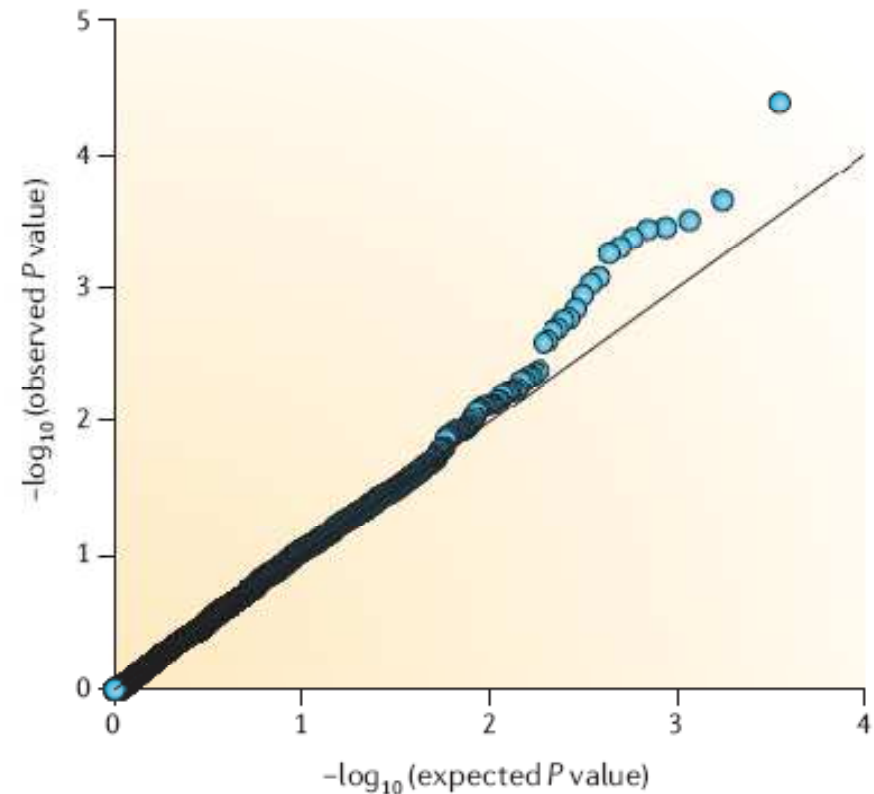
- Deviations from HWE can be due to inbreeding, population stratification or selection.
- Researchers have tested for HWE primarily as a data quality check and have discarded loci that, for example, deviate from HWE among controls at significance level $\alpha = 10^{-3}$ or 10^{-4} .
- Deviations from HWE can also be a symptom of disease association.
- So the possibility that a deviation from HWE is due to a deletion polymorphism or a segmental duplication that could be important in disease causation should certainly be considered before simply discarding loci...

Hardy-Weinberg equilibrium testing

- Testing for deviations from HWE can be carried out using a Pearson goodness-of-fit test, often known simply as ‘the χ^2 test’ because the test statistic has approximately a χ^2 null distribution.
- There are many different χ^2 tests. The Pearson test is easy to compute, but the χ^2 approximation can be poor when there are low genotype counts, in which case it is better to use a Fisher exact test.
- Fisher exact test does not rely on the χ^2 approximation.
- The open-source data-analysis software R has an *R genetics package* that implements both Pearson and Fisher tests of HWE

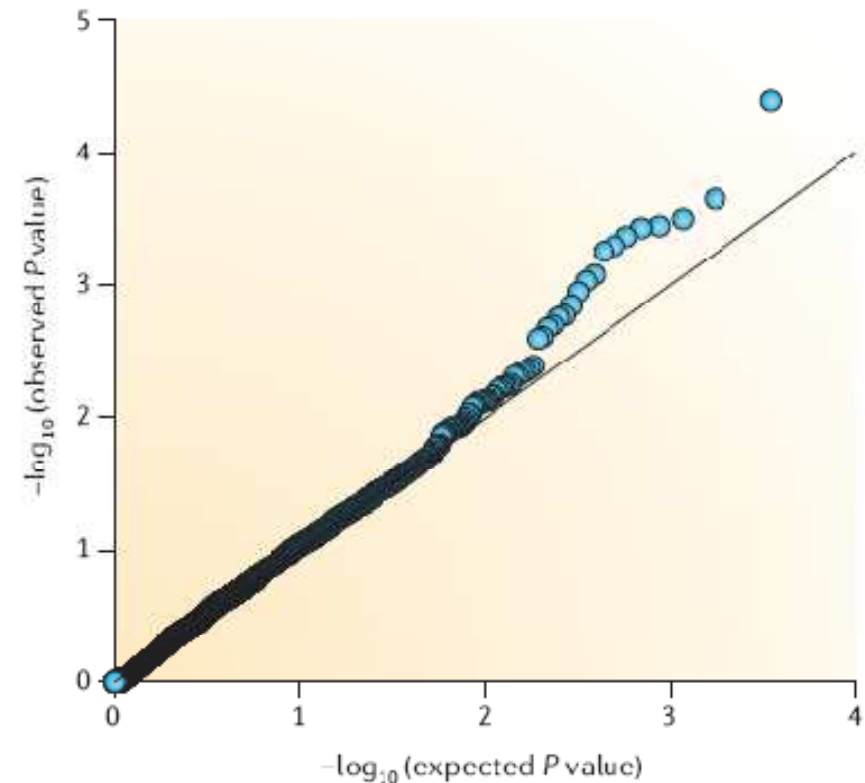
Hardy-Weinberg equilibrium interpretation of test results

- A useful tool for interpreting the results of HWE and other tests on many SNPs is the log quantile–quantile (QQ) p -value plot: the negative logarithm of the i -th smallest p -value is plotted against $-\log(i / (L + 1))$, where L is the number of SNPs.
- Deviations from the $y = x$ line correspond to loci that deviate from the null hypothesis.



Hardy-Weinberg equilibrium interpretation of test results

- The close adherence of p -values to the black line over most of the range is encouraging as it implies that there are few systematic sources of spurious association.
- The plot is suggestive of multiple weak associations, but the deviation of observed small p -values from the null line is unlikely to be sufficient to reach a reasonable criterion of significance.



2.c Missing genotype data

Introduction

- For single-SNP analyses
- , if a few genotypes are missing there is not much problem.
- For multipoint SNP analyses, missing data can be more problematic because many individuals might have one or more missing genotypes. One convenient solution is data imputation
- Data imputation involves replacing missing genotypes with predicted values that are based on the observed genotypes at neighbouring SNPs.
- For tightly linked markers data imputation can be reliable, can simplify analyses and allows better use of the observed data.
- For not tightly linked markers?

Introduction

- Imputation methods either seek a best prediction of a missing genotype, such as a
 - maximum-likelihood estimate (single imputation), or
 - randomly select it from a probability distribution (multiple imputations).
- The advantage of the latter approach is that repetitions of the random selection can allow averaging of results or investigation of the effects of the imputation on resulting analyses.
- Beware of settings in which cases are collected differently from controls. These can lead to differential rates of missingness even if genotyping is carried out blind to case-control status.
 - One way to check differential missingness rates is to code all observed genotypes as 1 and unobserved genotypes as 0 and to test for association of this variable with case-control status ...

2.d Haplotype and genotype data

Introduction

- Underlying an individual's genotypes at multiple tightly linked SNPs are the two haplotypes, each containing alleles from one parent.
- Analyses based on phased haplotype data rather than unphased genotypes may be *quite powerful*...

M1	1		1		2		2
DSL	D		d		d		d
M2	1		2		1		2

Test 1 vs. 2 for M1:

D + d vs. d

Test 1 vs. 2 for M2:

D + d vs. d

Test haplotype H1 vs. all others:

D vs. d

- If DSL located at a marker, haplotype testing can be *less powerful*

Inferring haplotypes

- Direct, laboratory-based haplotyping or typing further family members to infer the unknown phase are expensive ways to obtain haplotypes. Fortunately, there are statistical methods for inferring haplotypes and population haplotype frequencies from the genotypes of unrelated individuals.
- These methods, and the software that implements them, rely on the fact that in regions of low recombination relatively few of the possible haplotypes will actually be observed in any population.
- These programs generally perform well, given high SNP density and not too much missing data.

Inferring haplotypes

- Software:
 - **SNPHAP** is simple and fast, whereas **PHASE** tends to be more accurate but comes at greater computational cost.
 - **FASTPHASE** is nearly as accurate as PHASE but much faster.
- Whatever software is used, remember that true haplotypes are more informative than genotypes.
- Inferred haplotypes are typically less informative because of uncertain phasing.
 - The information loss that arises from phasing is small when linkage disequilibrium (LD) is strong.

Measures of LD

- LD will remain crucial to the design of association studies until whole-genome resequencing becomes routinely available. Currently, few of the more than 10 million common human polymorphisms are typed in any given study.
- If a causal polymorphism is not genotyped, we can still hope to detect its effects through LD with polymorphisms that are typed (key principle behind doing genetic association analysis ...).
- Hence, to assess the power of a study design to achieve this, we need to measure LD.

Measures of LD: D'

- LD is a non-quantitative phenomenon: there is no natural scale for measuring it.
- Among the measures that have been proposed for two-locus haplotype data, the two most important are D' (Lewontin's D prime) and r^2 (the square correlation coefficient between the two loci under study).
- The measure D is defined as the difference between the observed and expected (under the null hypothesis of independence) proportion of haplotypes bearing specific alleles at two loci: $p_{AB} - p_A p_B$

	A	a
B	p_{AB}	p_{aB}
b	p_{Ab}	p_{ab}

- D' is D/D_{\max}

Properties for D'

- D' is sensitive to even a few recombinations between the loci
- A disadvantage of D' is that it can be large (indicating high LD) even when one allele is very rare, which is usually of little practical interest.
- D' is inflated in small samples; the degree of bias will be greater for SNPs with rare alleles.
- So, the interpretation of values of $D' < 1$ is problematic, and values are difficult to compare between different samples because of the dependence on sample size.

Measures of LD: r^2

- r^2 is defined as

$$r^2 = \frac{D^2}{P_A P_a P_B P_b}$$

Properties for r^2

- In contrast to D' , r^2 is highly dependent upon allele frequency, and can be difficult to interpret when loci differ in their allele frequencies
- However, r^2 has desirable sampling properties, is directly related to the amount of information provided by one locus about the other, and is particularly useful in evolutionary and population genetics applications.
- Specifically, sample size must be increased by a factor of $1/r^2$ to detect an unmeasured variant, compared with the sample size for testing the variant itself.

(Jorgenson and Witte 2006)

1.e SNP tagging

Introduction

- Tagging refers to methods to select a minimal number of SNPs that retain as much as possible of the genetic variation of the full SNP set.
- Simple pairwise methods discard one (preferably that with most missing values) of every pair of SNPs with, say, $r^2 > 0.9$.
- More sophisticated methods can be more efficient, but the most efficient tagging strategy will depend on the statistical analysis to be used afterwards.
- In practice, tagging is only effective in capturing common variants.

Two good reasons for tagging

- The first principal use for tagging is to select a ‘good’ subset of SNPs to be typed in all the study individuals from an extensive SNP set that has been typed in just a few individuals.
 - Until recently, this was frequently a laborious step in study design, but the International HapMap Project and related projects now allow selection of tag SNPs on the basis of publicly available data.
 - However, the population that underlies a particular study will typically differ from the populations for which public data are available, and a set of tag SNPs that have been selected in one population might perform poorly in another.
 - Nevertheless, recent studies indicate that tag SNPs often transfer well across populations

Two good reasons for tagging

- The second use for tagging is to select for analysis a subset of SNPs that have already been typed in all the study individuals.
- Although it is undesirable to discard available information, the amount of information lost might be small (at least, that is what is aimed for when applying SNP tagging algorithms).
- Reducing the SNP set can simplify analyses and lead to more statistical power by reducing the degrees of freedom (df) of a test.

3 Tests of association: single SNP

Introduction

- Population association studies compare unrelated individuals, but 'unrelated' actually means that relationships are unknown and presumed to be distant.
- Therefore, we cannot trace transmissions of phenotype over generations and must rely on correlations of current phenotype with current marker alleles.
- Such a correlation might be generated (but is not necessarily generated) by one or more groups of cases that share a relatively recent common ancestor at a causal locus.

A toy example

	AA	AB	BB	total
(A) Genotype counts				
Case	a = 10	b = 190	c = 800	a+b+c = 1000
Control	d = 3	e = 100	f = 900	d+e+f = 1003

	A	B	total
(B) Allele counts			
Case	$x_{11} = 2a+b = 210$	$x_{12} = b+2c = 1790$	$2(a+b+c) = 2000$
Control	$x_{21} = 2d+e = 106$	$x_{22} = d+2f = 1900$	$2(d+e+f) = 2006$

	AA+AB	BB	AA	AB+BB	total
(C) Two ways of grouping heterozygotes with homozygotes					
Case	a+b = 200	c = 800	a = 10	b+c = 990	a+b+c = 1000
Control	d+e = 103	f = 900	c = 3	d+e = 1000	d+e+f = 1003

There are 1000 case samples and 1003 control samples, whose genotype distribution is shown in the table (A); the number of A and B allele counts is in (B). The genotype counts in (C) are converted from (A) by combining AB with either AA or BB. Note that the total counts in (B) doubles the counts in (A), and the two tables in (C) correspond to the dominant and recessive models if allele A is considered as the risk allele.

(Li 2007)

A toy example

- A Pearson's test is a summary of discrepancy between the observed (O) and expected (E) genotype/allele count:

$$\chi^2 = \sum_{i=1} \frac{(O_i - E_i)^2}{E_i}$$

- For any χ^2 distributed test statistic with df degrees of freedom, one can decompose it to two χ^2 distributed test statistics with df 1 and df 2 degrees of freedom and their sum df 1 + df 2 is equal to df.
- For example, the test statistic in the genotype based test (GBT) can be decomposed to two χ^2 distributed values each with one degree of freedom.
- One of them is the test statistic in a commonly used test called Conchran–Armitage test (CAT).

A toy example

- CAT tests whether $\log(r)$, where r is the (number of cases)/(number of cases + number of controls) ratio, changes linearly with the AA, AB, BB genotype with a non-zero slope.
- Note that since AB is positioned between AA and BB genotype, the genotype is not just a categorical variable, but an ordered categorical variable.
- Also note that although CAT is genotype based, its value is closer to the allele-based ABT test statistic.

A toy example: testing

```
gc <- c(10, 190, 800, 3, 100, 900)
ac <- c(2*gc[1]+gc[2], gc[2]+2*gc[3], 2*gc[4]+gc[5], gc[5]+2*gc[6])
gc1 <- c(gc[1]+gc[2], gc[3], gc[4]+gc[5], gc[6])
gc2 <- c(gc[1],gc[2]+gc[3], gc[4], gc[5]+gc[6])
pvg <- chisq.test(matrix(gc, ncol=3, byrow=T), corr=FALSE)$p.value
pva <- chisq.test(matrix(ac, ncol=2, byrow=T), corr=FALSE)$p.value
pvg1 <- chisq.test(matrix(gc1, ncol=2, byrow=T), corr=FALSE)$p.value
pvg2 <- chisq.test(matrix(gc2, ncol=2, byrow=T), corr=FALSE)$p.value
pvb <- min(pvg1, pvg2)

print(c(pvg, pva, pvb)) # 6.918239e-09 9.150309e-10 1.224003e-09
pvg.f <- fisher.test(matrix(gc, ncol=3, byrow=T))$p.value
pva.f <- fisher.test(matrix(ac, ncol=2, byrow=T))$p.value
pvg1.f <- fisher.test(matrix(gc1, ncol=2, byrow=T))$p.value
pvg2.f <- fisher.test(matrix(gc2, ncol=2, byrow=T))$p.value
pvb.f <- min(pvg1.f, pvg2.f)
print(c(pvg.f, pva.f, pvb.f)) # 2.412721e-09 8.047005e-10 1.132535e-09

pvcat <- prop.trend.test(gc[1:3], gc[1:3]+gc[4:6], score=c(0, 0.5, 1))$p.value
print(c(pvcat) ) # 9.820062e-10

gc <- gc*2
... # repeat the tests
print(c(pvg, pva, pvb)) # 4.786203e-17 4.716312e-18 8.379499e-18
print(c(pvg.f, pva.f, pvb.f)) # 1.231881e-17 3.485271e-18 6.810263e-18
print(c(pvcat) ) # 5.422705e-18
```

A toy example: testing

- What is the effect of choosing a different genetic model?
- What is the effect of choosing a genotype test versus an allelic test?
- Are allelic tests always applicable?
- When do you expect the largest differences between Pearson's chi-square and Fisher's exact test?
- What is the effect of doubling the sample size on these tests?
- How can you protect yourself against uncertain disease models?

A toy example: estimation

```
ci.or <- function(counts, alpha){ # alpha=0.05 corresponds to 95%CI
  f <- qnorm(1- alpha/2)         # if alpha=0.05, f=1.96
  or <- counts[1]*counts[4]/(counts[2]*counts[3])
  sq <- sqrt(1/counts[1]+1/counts[2]+1/counts[3]+1/counts[4])
  upper <- exp( log(or) + f*sq)
  lower <- exp( log(or) - f*sq)
  res <- c(lower, or, upper)
  res
}

print( ci.or(ac, 0.05))          # 1.650411 2.102878 2.679390
print( ci.or(ac, 0.01))          # 1.529428 2.102878 2.891339

ac <- ac*2                       # double the sample size
print( ci.or(ac, 0.05))          # 1.771784 2.102878 2.495842
print( ci.or(ac, 0.01))          # 1.678927 2.102878 2.633882
```

A toy example: estimation

- Will all packages give you the same output when estimating odds ratios with confidence intervals, assuming the data and the significance level are the same?
- What is the effect of decreasing the significance level?
- What is the effect of doubling the sample size?

Example R code to perform small-scale analyses using GENETICS

```
library(DGCgenetics)
library(dgc.genetics)
casecon <- read.table("casecondata.txt",header=T)
casecon[1:2,]
attach(casecon)
pedigree
case <- affected-1
case
g1 <- genotype(loc1_1,loc1_2)
g1 <- genotype(loc2_1,loc2_2)
g1 <- genotype(loc3_1,loc3_2)
g1 <- genotype(loc1_1,loc1_2)
g2 <- genotype(loc2_1,loc2_2)
g3 <- genotype(loc3_1,loc3_2)
g4 <- genotype(loc4_1,loc4_2)
g1
```

```
table(g1,case)
chisq.test(g1,case)
allele.table(g1,case)
gcontrasts(g1) <- "genotype"
names(casecon)
help(gcontrasts)
logit(case~g1)
anova(logit(case~g1))
1-pchisq(18.49,2)
gcontrasts(g1) <- "genotype"
gcontrasts(g3) <- "genotype"
logit(case~g1+g3)
anova(logit(case~g1+g3))
gcontrasts(g1) <- "genotype"
gcontrasts(g3) <- "additive"
logit(case~g1+g3)
anova(logit(case~g1+g3))
detach(casecon)
```

This is in fact already a multiple SNP analysis
But you can see how easy it is within a
regression framework

Example R code to perform small-scale analyses using SNPassoc

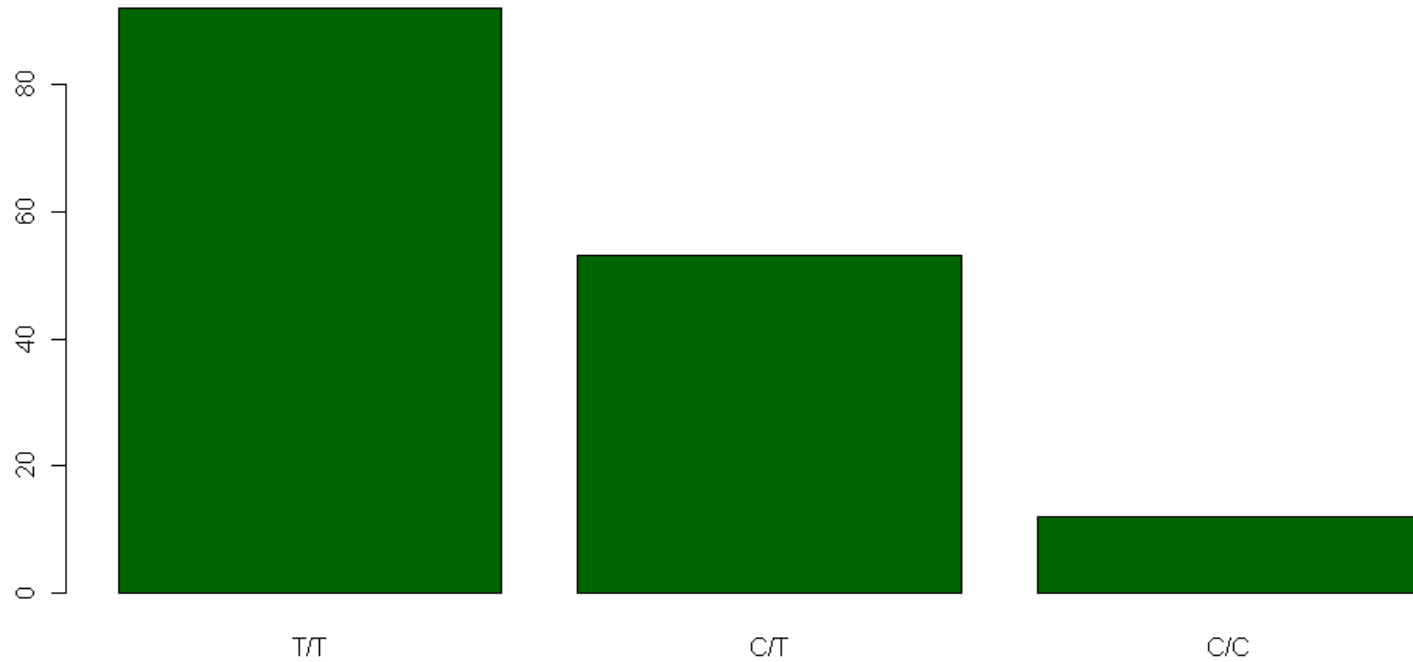
```
#Let's load library SNPassoc
library(SNPassoc)
#get the data example:
#both data.frames SNPs and SNPs.info.pos are loaded typing data(SNPs)
data(SNPs)
#look at the data (only first four SNPs)
SNPs[1:10,1:9]
table(SNPs[,2])
mySNP<-snp(SNPs$snp10001,sep="")
mySNP
summary(mySNP)
```

snp10001

	frequency	percentage
T	237	75.48
C	77	24.52

	frequency	percentage
T/T	92	58.60
C/T	53	33.76
C/C	12	7.64

HWE (pvalue): 0.281639



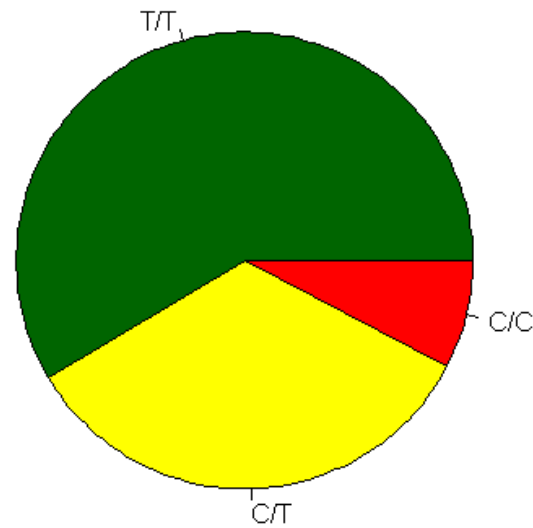
```
plot(mySNP,label="snp10001",col="darkgreen")
```

snp10001

	frequency	percentage
T	237	75.48
C	77	24.52

	frequency	percentage
T/T	92	58.60
C/T	53	33.76
C/C	12	7.64

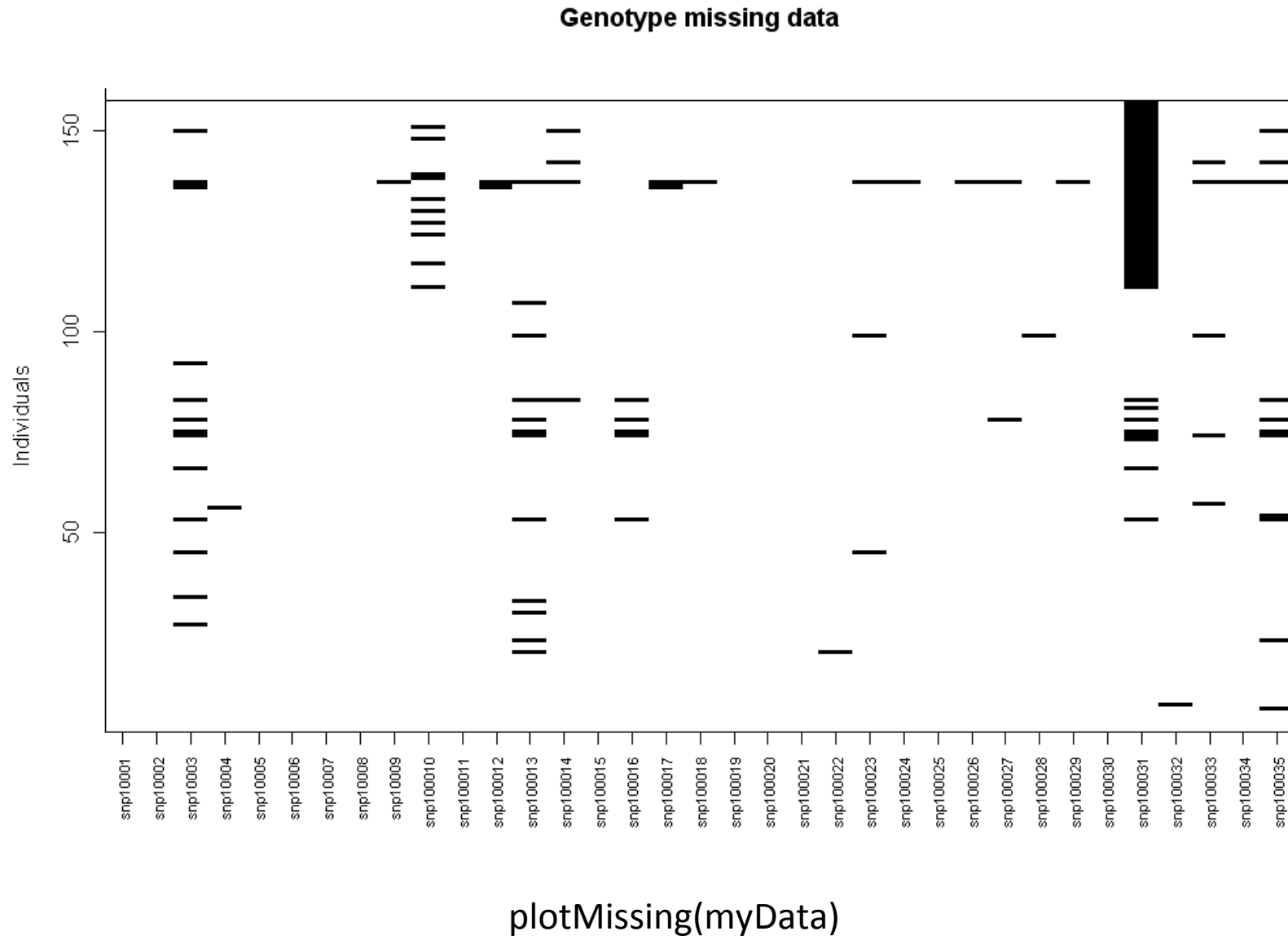
HWE (pvalue): 0.281639



```
plot(mySNP,type=pie,label="snp10001",col=c("darkgreen","yellow","red"))
```

Example R code to perform small-scale analyses using SNPassoc

```
reorder(mySNP,ref="minor")
gg<-
c("het","hom1","hom1","hom1","hom1","hom1","het","het","het","hom1","hom2","hom
1","hom2")
snp(gg,name.genotypes=c("hom1","het","hom2"))
myData<-setupSNP(data=SNPs,colSNPs=6:40,sep="")
myData.o<-setupSNP(SNPs, colSNPs=6:40, sort=TRUE,info=SNPs.info.pos, sep="")
labels(myData)
summary(myData)
plot(myData,which=20)
```



Example R code to perform small-scale analyses using SNPassoc

```
res<-tableHWE(myData)
res
res<- tableHWE(myData,strata=myData$sex)
res
```

What is the difference between the two previous commands?
Why is the latter analysis important?

Example R code to perform GWA using SNPassoc

```
data(HapMap)
```

```
> HapMap[1:4,1:9]
```

```
      id group rs10399749 rs11260616 rs4648633 rs6659552 rs7550396 rs12239794
rs6688969
1 NA06985 CEU      CC      AA      TT      GG      GG      GG      CC
2 NA06993 CEU      CC      AT      CT      CG      GG      GG      CT
3 NA06994 CEU      CC      AA      TT      CG      GG      GG      CT
4 NA07000 CEU      CC      AT      TT      GG      GG      <NA>  CC
```

```
myDat.HapMap<-setupSNP(HapMap, colSNPs=3:9307, sort =
TRUE,info=HapMap.SNPs.pos, sep="")
```

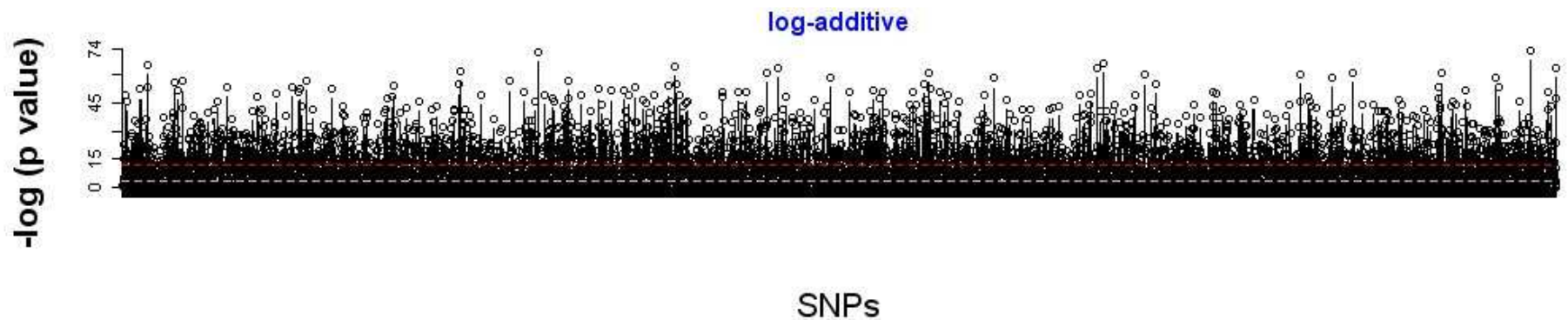
```
> HapMap.SNPs.pos[1:3,]
```

```
      snp chromosome position
1 rs10399749      chr1  45162
2 rs11260616      chr1 1794167
3 rs4648633       chr1 2352864
```

Example R code to perform GWA using SNPAssoc

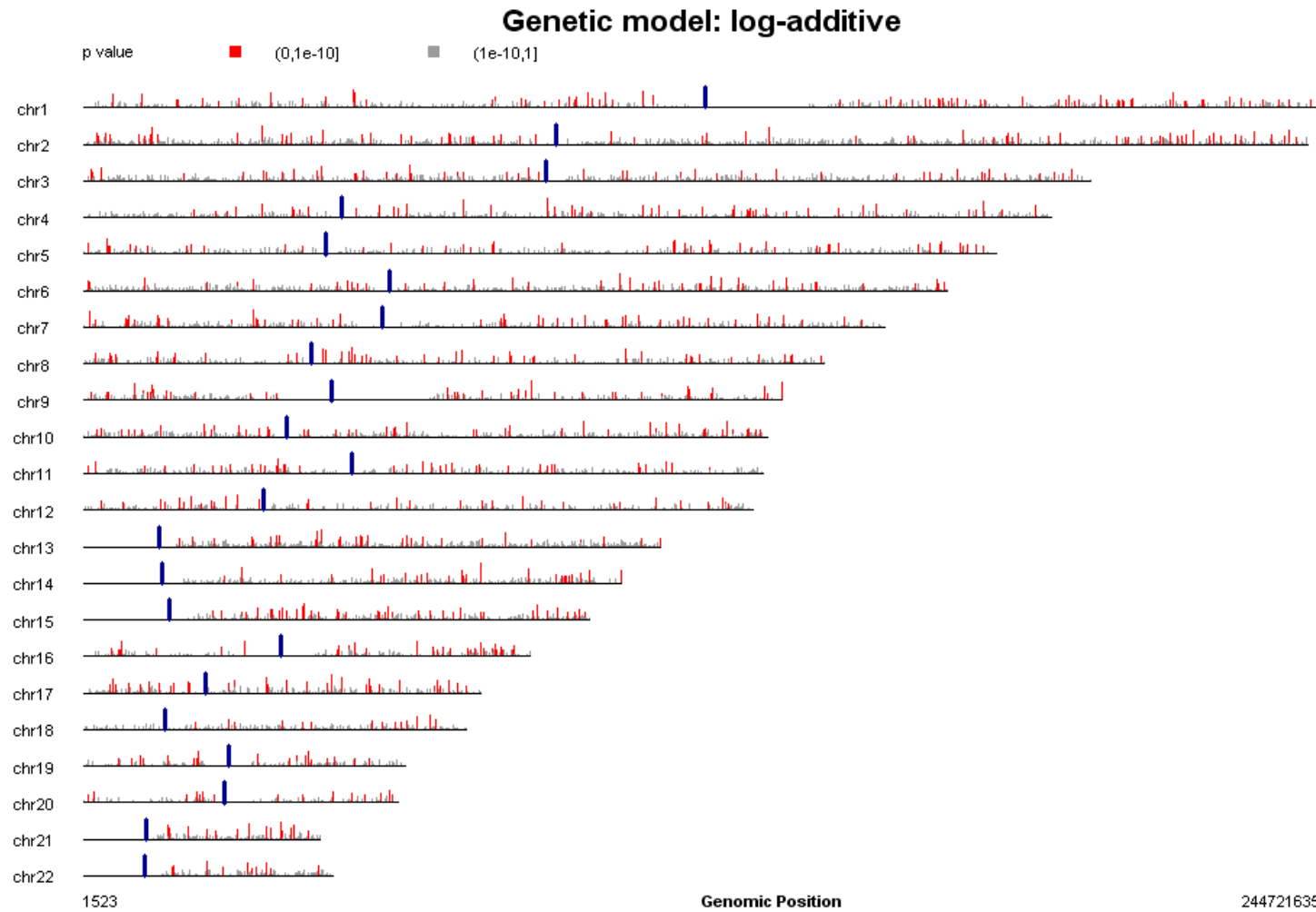
```
resHapMap<-WGassociation(group, data=myDat.HapMap, model="log-add")
```

```
plot(resHapMap, whole=FALSE, print.label.SNPs = FALSE)
```



```
> summary(resHapMap)
```

SNPs (n)	Genot error (%)	Monomorphic (%)	Significant* (n)	(%)
chr1 796	3.8	18.6	163	20.5
chr2 789	4.2	13.9	161	20.4
chr3 648	5.2	13.0	132	20.4



`plot(resHapMap, whole=TRUE, print.label.SNPs = FALSE)`

Example R code to perform GWA using SNPassoc

```
resHapMap.scan<-scanWGassociation(group, data=myDat.HapMap, model="log-add")
resHapMap.perm<-scanWGassociation(group, data=myDat.HapMap,model="log-add",
nperm=1000)
res.perm<- permTest(resHapMap.perm)
```

- Check out the SNPassoc manual (supporting document to R package) to read more about the analytical methods used

Example R code to perform GWA using SNPAssoc

```
> print(resHapMap.scan[1:5,])
      comments log-additive
rs10399749 Monomorphic -
rs11260616 -      0.34480
rs4648633  -      0.00000
rs6659552  -      0.00000
rs7550396  -      0.31731
> print(resHapMap.perm[1:5,])
      comments log-additive
rs10399749 Monomorphic -
rs11260616 -      0.34480
rs4648633  -      0.00000
rs6659552  -      0.00000
rs7550396  -      0.31731
```

```
perms <- attr(resHapMap.perm, "pvalPerm") #what does this object contain?
```

Example R code to perform GWA using SNPassoc

```
> print(res.perm)
```

Permutation test analysis (95% confidence level)

Number of SNPs analyzed: 9305

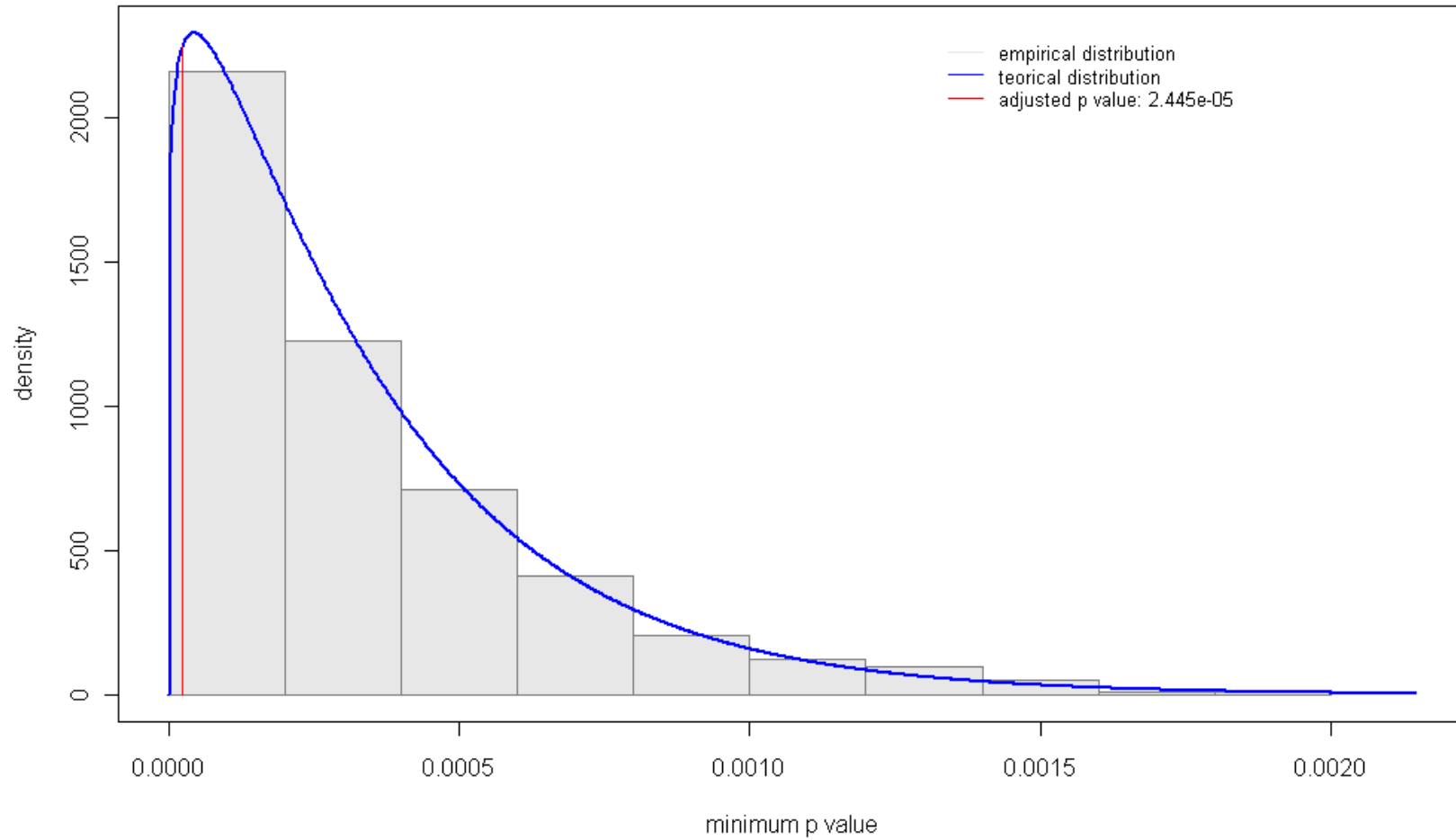
Number of valid SNPs (e.g., non-Monomorphic and passing calling rate): 7320

P value after Bonferroni correction: 6.83e-06

P values based on permutation procedure:

P value from empirical distribution of minimum p values: 2.883e-05

P value assuming a Beta distribution for minimum p values: 2.445e-05



plot(res.perm)

Example R code to perform GWA using SNPassoc

```
res.perm.rtp<- permTest(resHapMap.perm,method="rtp",K=20)  
> print(res.perm.rtp)
```

Permutation test analysis (95% confidence level)

Number of SNPs analyzed: 9305

Number of valid SNPs (e.g., non-Monomorphic and passing calling rate):
7320

P value after Bonferroni correction: 6.83e-06

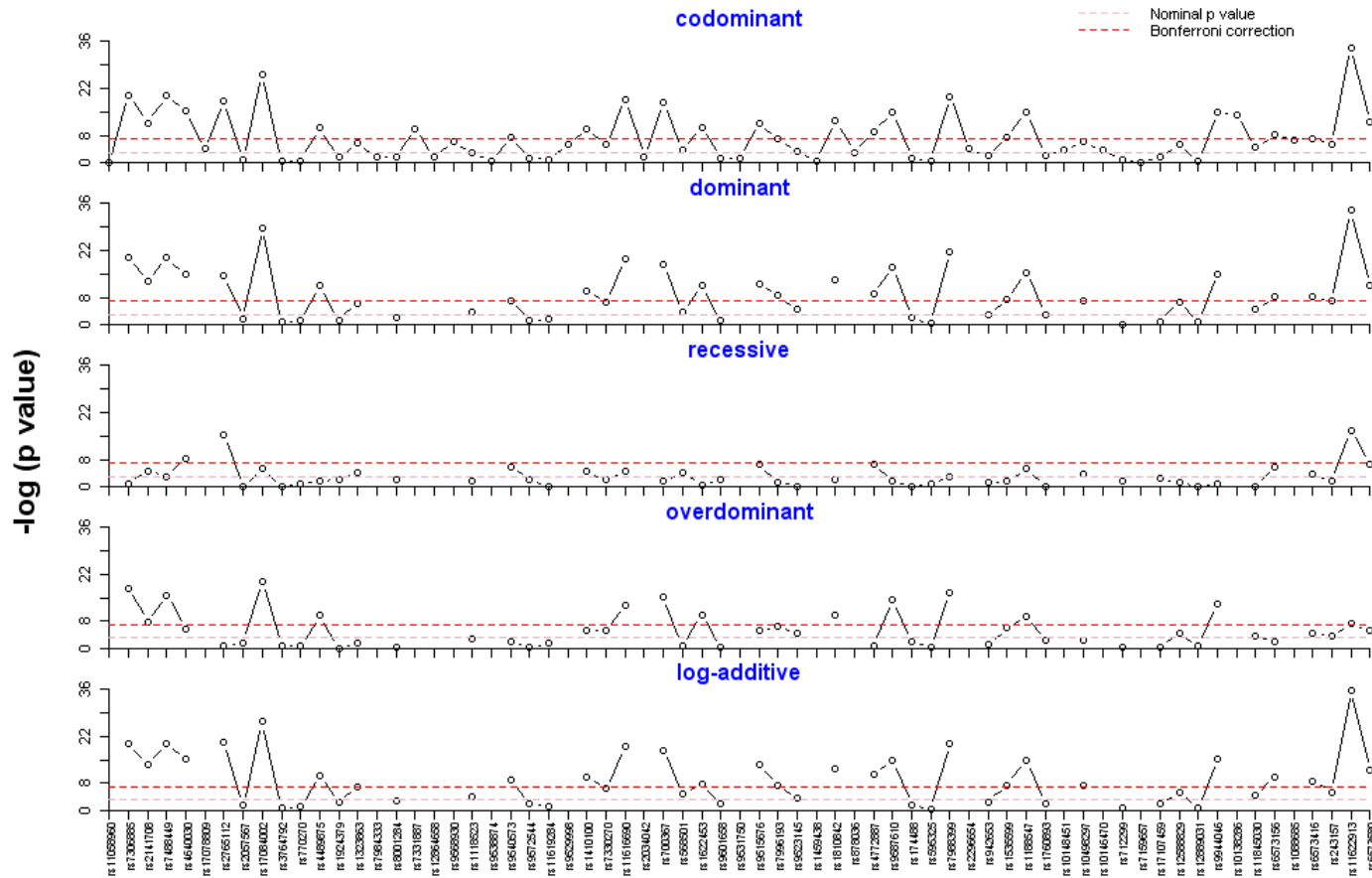
Rank truncated product of the K=20 most significant p-values:

Product of K p-values (-log scale): 947.2055

Significance: <0.001

Example R code to perform a variety of medium/large-scale analyses using SNPAssoc

```
getSignificantSNPs(resHapMap,chromosome=5)
association(casco~snp(snp10001,sep=""), data=SNPs)
myData<-setupSNP(data=SNPs,colSNPs=6:40,sep="")
association(casco~snp10001, data=myData)
association(casco~snp10001, data=myData, model=c("cod","log"))
association(casco~sex+snp10001+blood.pre, data=myData)
association(casco~snp10001+blood.pre+strata(sex), data=myData)
association(casco~snp10001+blood.pre, data=myData,subset=sex=="Male")
association(log(protein)~snp100029+blood.pre+strata(sex), data=myData)
ans<-association(log(protein)~snp10001*sex+blood.pre,
data=myData,model="codominant")
print(ans,dig=2)
ans<-association(log(protein)~snp10001*factor(recessive(snp100019))+blood.pre,
data=myData, model="codominant")
print(ans,dig=2)
```



```
sigSNPs<-getSignificantSNPs(resHapMap,chromosome=5,sig=5e-8)$column
myDat2<-setupSNP(HapMap, colSNPs=sigSNPs, sep="")
resHapMap2<-WGassociation(group~1, data=myDat2)
plot(resHapMap2,cex=0.8)
```

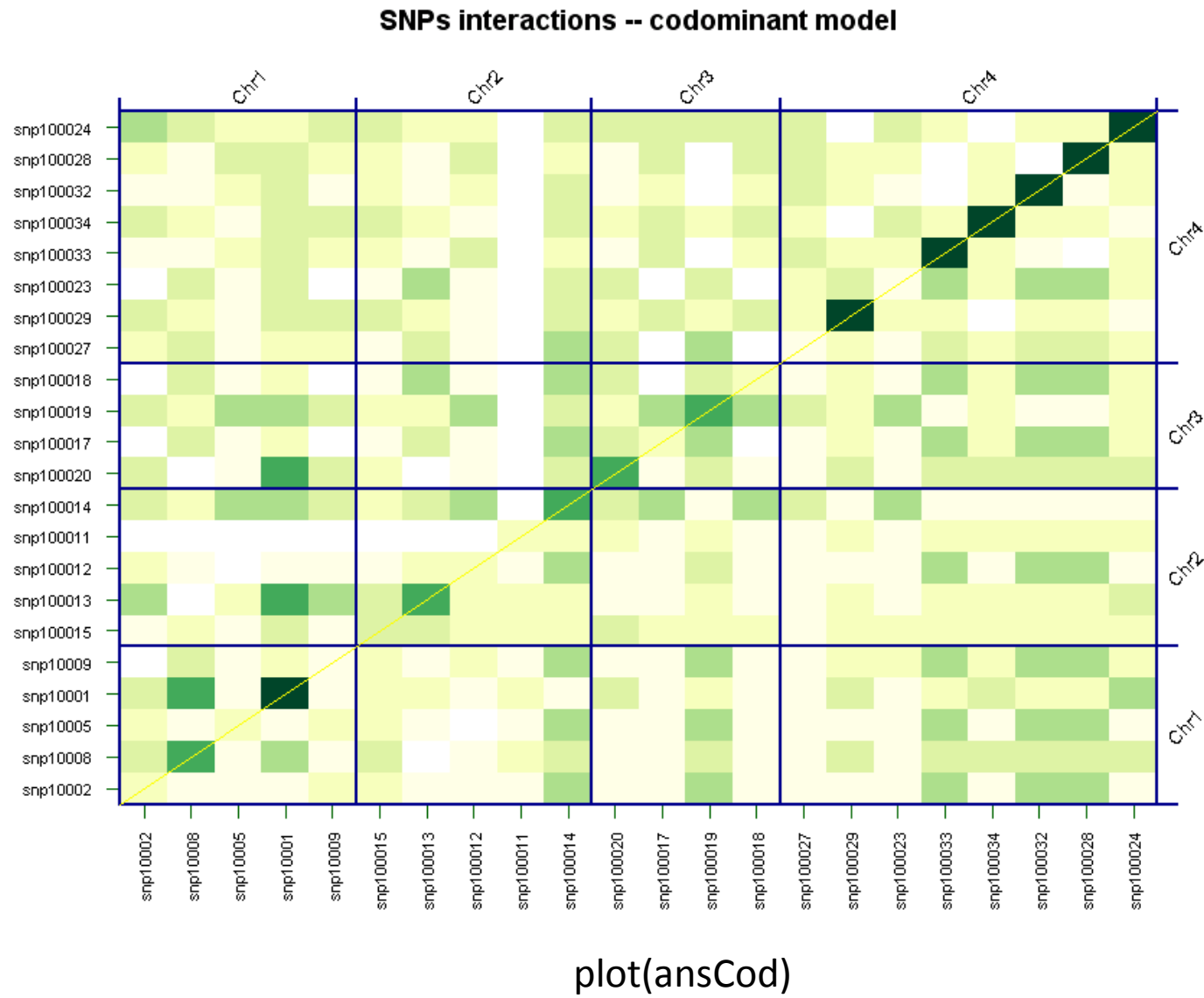

4 Tests of association: multiple SNPs

Introduction

- Choices to be made:
 - Enter multiple markers in one model
 - Analyze the markers as independent contributors (see earlier example R code)
 - Analyze the markers as potentially interacting (see Chapter 9)
 - Construct haplotypes from multiple tightly linked markers and analyze accordingly
- All these analyses are easily performed in a “regression” context
 - In particular, for case / control data, logistic regression is used, where disease status is regressed on genetic predictors

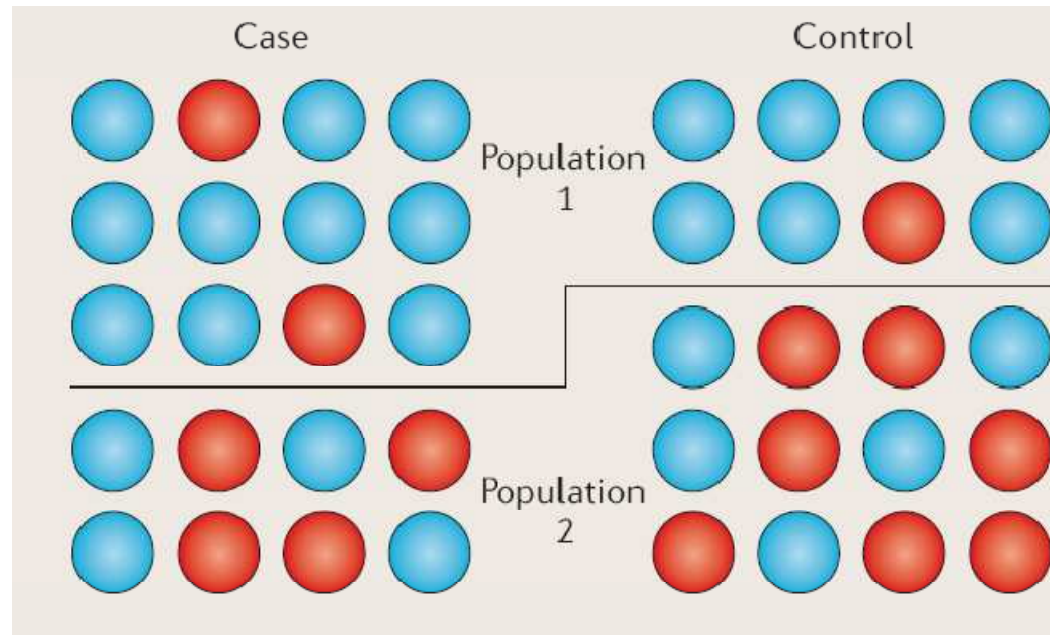
Example R code using SNPassoc

```
datSNP<-setupSNP(SNPs,6:40,sep="")
tag.SNPs<-c("snp100019", "snp10001", "snp100029")
geno<-make.geno(datSNP,tag.SNPs)
mod<-
haplo.glm(log(protein)~geno,data=SNPs,family=gaussian,locus.label=tag.SNPs,allele.lev=at
tributes(geno)$unique.alleles,
control = haplo.glm.control(haplo.freq.min=0.05))
mod
intervals(mod)
ansCod<-interactionPval(log(protein)~sex, data=myData.o,model="codominant")
```



5 Dealing with population stratification

5.a Spurious associations



- Methods to deal with spurious associations generated by population structure generally require a number (preferably >100) of widely spaced null SNPs that have been genotyped in cases and controls in addition to the candidate SNPs.

5.b Genomic Control

- In Genomic Control (GC), a 1-df association test statistic (usually, CAT) is computed at each of the null SNPs, and a parameter λ is calculated as the empirical median divided by its expectation under the chi-squared 1-df distribution.
- Then the association test is applied at the candidate SNPs, and if $\lambda > 1$ the test statistics are divided by λ .
- There is an analogous procedure for a general (2 df) test; The method can also be applied to other testing approaches.
- The motivation for GC is that, as we expect few if any of the null SNPs to be associated with the phenotype, a value of $\lambda > 1$ is likely to be due to the effect of population stratification, and dividing by λ cancels this effect for the candidate SNPs.
- GC performs well under many scenarios, but can be conservative in extreme settings (and anti-conservative if insufficient null SNPs are used).

5.c Structured Association methods

- Structured association (SA) approaches are based on the idea of attributing the genomes of study individuals to hypothetical subpopulations, and testing for association that is conditional on this subpopulation allocation.
- These approaches are computationally demanding, and because the notion of subpopulation is a theoretical construct that only imperfectly reflects reality, the question of the correct number of subpopulations can never be fully resolved....

5.d Other approaches to handle the effects of population substructure

Include extra covariates in regression models used for association modeling/testing

- Null SNPs can mitigate the effects of population structure when included as covariates in regression analyses.
- Like GC, this approach does not explicitly model the population structure and is computationally fast, but it is much more flexible than GC because epistatic and covariate effects can be included in the regression model.
- Empirically, the logistic regression approaches show greater power than GC, but their type-1 error rate must be determined through simulation.
- Simulations can be quite intensive! How many replicates are sufficient?

Principal components analysis

- When many null markers are available, principal components analysis provides a fast and effective way to diagnose population structure.
- In European data, the first 2 principal components nicely reflect the N-S and E-W axes

Unrelateds are “distantly” related

- Alternatively, a mixed-model approach that involves estimated kinship, with or without an explicit subpopulation effect, has recently been found to outperform GC in many settings.
- Given large numbers of null SNPs, it becomes possible to make precise statements about the (distant) relatedness of individuals in a study so that in theory it should be possible to provide a complete solution to the problem of population stratification.

6 Multiple testing

6.a General setting

Introduction

- Multiple testing is a thorny issue, the bane of statistical genetics.
 - The problem is not really the number of tests that are carried out: even if a researcher only tests one SNP for one phenotype, if many other researchers do the same and the nominally significant associations are reported, there will be a problem of false positives.
- The genome is large and includes many polymorphic variants and many possible disease models. Therefore, any given variant (or set of variants) is highly unlikely, *a priori*, to be causally associated with any given phenotype under the assumed model.
- So strong evidence is required to overcome the appropriate scepticism about an association.

6.b Controlling the overall type I error

Frequentist paradigm

- The frequentist paradigm of controlling the overall type-1 error rate sets a significance level α (often 5%), and all the tests that the investigator plans to conduct should together generate no more than probability α of a false positive.
- In complex study designs, which involve, for example, multiple stages and interim analyses, this can be difficult to implement, in part because it was the analysis that was planned by the investigator that matters, not only the analyses that were actually conducted.

Frequentist paradigm

- In simple settings the frequentist approach gives a practical prescription:
 - if n SNPs are tested and the tests are approximately independent, the appropriate per-SNP significance level α' should satisfy

$$\alpha = 1 - (1 - \alpha')^n,$$

which leads to the Bonferroni correction $\alpha' \approx \alpha / n$.

- For example, to achieve $\alpha = 5\%$ over 1 million independent tests means that we must set $\alpha' = 5 \times 10^{-8}$. However, the *effective number* of independent tests in a genome-wide analysis depends on many factors, including sample size and the test that is carried out.

When markers (and hence tests) are tightly linked

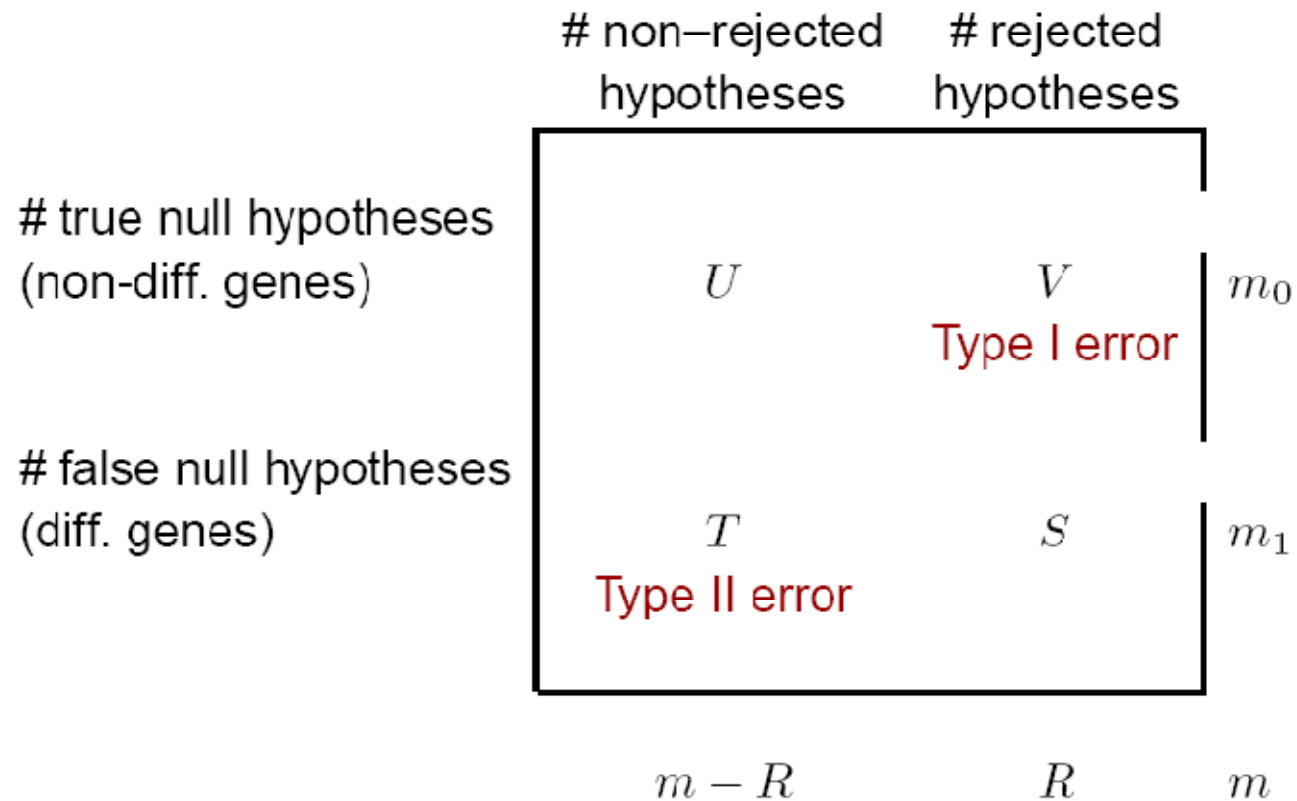
- For tightly linked SNPs, the Bonferroni correction is conservative.
- A practical alternative is to approximate the type-I error rate using a permutation procedure.
 - Here, the genotype data are retained but the phenotype labels are randomized over individuals to generate a data set that has the observed LD structure but that satisfies the null hypothesis of no association with phenotype.
 - By analysing many such data sets, the false-positive rate can be approximated.
 - The method is conceptually simple but can be computationally demanding, particularly as it is specific to a particular data set and the whole procedure has to be repeated if other data are considered.

The 5% magic percentage

- Although the 5% global error rate is widely used in science, it is inappropriately conservative for large-scale SNP-association studies:
 - Most researchers would accept a higher risk of a false positive in return for greater power.
- There is no “rule” saying that the 5% value cannot be relaxed, but another approach is to monitor the false discovery rate (FDR) instead
- The FDR refers to the *proportion of false positive test results among all positives*.

FDR control

- In particular,



(Benjamini and Hochberg 1995: $FDR = E(Q)$; $Q = V/R$ when $R > 0$ and $Q = 0$ when $R = 0$)

FDR control

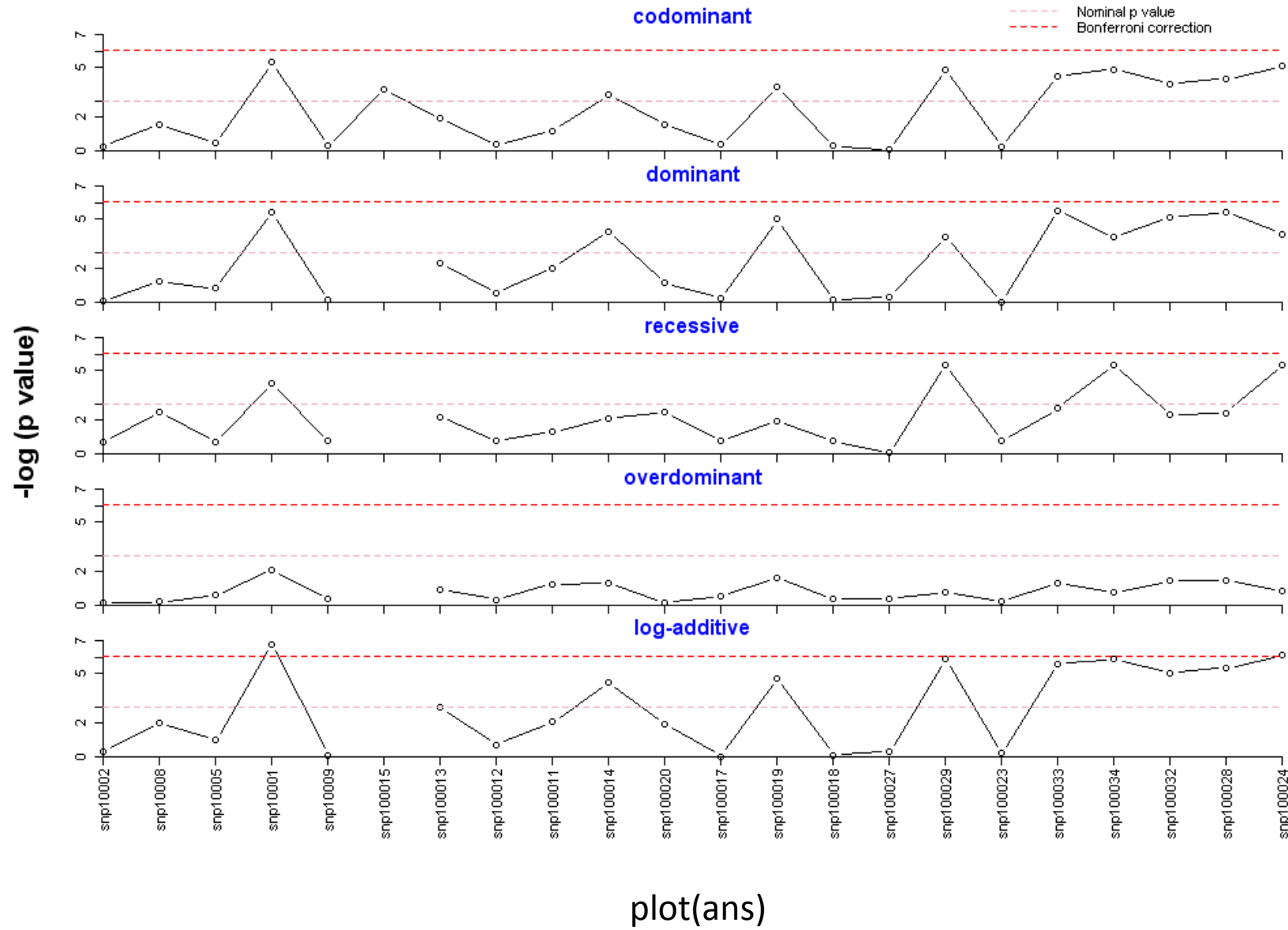
- FDR measures come in different shapes and flavor.
 - But under the null hypothesis of no association, p -values should be uniformly distributed between 0 and 1;
 - FDR methods typically consider the actual distribution as a mixture of outcomes under the null (uniform distribution of p -values) and alternative (P -value distribution skewed towards zero) hypotheses.
 - Assumptions about the alternative hypothesis might be required for the most powerful methods, but the simplest procedures avoid making these explicit assumptions.

Cautionary note

- The usual frequentist approach to multiple testing has a serious drawback in that researchers might be discouraged from carrying out additional analyses beyond single-SNP tests, even though these might reveal interesting associations, because all their analyses would then suffer a multiple-testing penalty.
- It is a matter of common sense that expensive and hard-won data should be investigated exhaustively for possible patterns of association.
- Although the frequentist paradigm is convenient in simple settings, strict adherence to it can be dangerous: true associations may be missed!
 - Under the Bayesian approach, there is no penalty for analysing data exhaustively because the prior probability of an association should not be affected by what tests the investigator chooses to carry out.

Example R code using SNPassoc

```
myData<-setupSNP(SNPs, colSNPs=6:40, sep="")
myData.o<-setupSNP(SNPs, colSNPs=6:40, sort=TRUE,info=SNPs.info.pos, sep="")
ans<-WGassociation(protein~1,data=myData.o)
library(Hmisc)
SNP<-pvalues(ans)
out<-latex(SNP,file="c:/temp/ans1.tex", where=""h",caption="Summary of case-control
study for SNPs data set.",center="centering", longtable=TRUE, na.blank=TRUE,
size="scriptsize", collabel.just=c("c"), lines.page=50,rownamesTexCmd="bfseries")
WGstats(ans,dig=5)
```



Example R code using SNPassoc

```
Bonferroni.sig(ans, model="log-add", alpha=0.05,include.all.SNPs=FALSE)
```

```
pvalAdd<-additive(resHapMap)
```

```
pval<-pval[!is.na(pval)]
```

```
library(qvalue)
```

```
qobj<-qvalue(pval)
```

```
max(qobj$qvalues[qobj$pvalues <= 0.001])
```

```
procs<-c("Bonferroni","Holm","Hochberg","SidakSS","SidakSD","BH","BY")
```

```
res2<-mt.rawp2adjp(rawp,procs)
```

```
mt.reject(cbind(res$rawp,res$adjp),seq(0,0.1,0.001))$r
```

7 Assessing the function of genetic variants

Criteria for assessing the functional significance of a variant

Criteria	Strong support for functional significance	Moderate support for functional significance	Evidence against functional significance
Nucleotide sequence	Variant disrupts a known functional or structural motif	Variant is a missense change or disrupts a putative functional motif; changes to protein structure might occur	Variant disrupts a non-coding region with no known functional or structural motif
Evolutionary conservation	Consistent evidence from multiple approaches for conservation across species and multigene families	Evidence for conservation across species or multigene families	Nucleotide or amino-acid residue not conserved
Population genetics	In the absence of laboratory error, strong deviations from expected population frequencies in cases and/or controls in a particular ethnicity	In the absence of laboratory error, moderate to small deviations from expected population frequencies in cases and/or controls; effects are not well characterized by ethnicity	Population genetics data indicates no deviations from expected proportions
Experimental evidence	Consistent effects from multiple lines of experimental evidence; effect in human context is established; effect in target tissue is known	Some (possibly inconsistent) evidence for function from experimental data; effect in human context or target tissue is unclear	Experimental evidence consistently indicates no functional effect
Exposures (for example, genotype-environment interaction studies)	Variant is known to affect the metabolism of the exposure in the relevant target tissue	Variant might affect metabolism of the exposure or one of its components; effect in target tissue might not be known	Variant does not affect metabolism of exposure of interest
Epidemiological evidence	Consistent and reproducible reports of moderate-to-large magnitude associations	Reports of association exist; replication studies are not available	Prior studies show no effect of variant

(Rebbeck et al 2004)

8 Proof of concept

Vol 447 | 7 June 2007 | doi:10.1038/nature05911

nature

ARTICLES

Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

The Wellcome Trust Case Control Consortium*

There is increasing evidence that genome-wide association (GWA) studies represent a powerful approach to the identification of genes involved in common human diseases. We describe a joint GWA study (using the Affymetrix GeneChip 500K Mapping Array Set) undertaken in the British population, which has examined ~2,000 individuals for each of 7 major diseases and a shared set of ~3,000 controls. Case-control comparisons identified 24 independent association signals at $P < 5 \times 10^{-7}$: 1 in bipolar disorder, 1 in coronary artery disease, 9 in Crohn's disease, 3 in rheumatoid arthritis, 7 in type 1 diabetes and 3 in type 2 diabetes. On the basis of prior findings and replication studies thus-far completed, almost all of these signals reflect genuine susceptibility effects. We observed association at many previously identified loci, and found compelling evidence that some loci confer risk for more than one of the diseases studied. Across all diseases, we identified a large number of further signals (including 58 loci with single-point P values between 10^{-5} and 5×10^{-7}) likely to yield additional susceptibility loci. The importance of appropriately large samples was confirmed by the modest effect sizes observed at most loci identified. This study thus represents a thorough validation of the GWA approach. It has also demonstrated that careful use of a shared control group represents a safe and effective approach to GWA analyses of multiple disease phenotypes; has generated a genome-wide genotype database for future studies of common diseases in the British population; and shown that, provided individuals with non-European ancestry are excluded, the extent of population stratification in the British population is generally modest. Our findings offer new avenues for exploring the pathophysiology of these important disorders. We anticipate that our data, results and software, which will be widely available to other investigators, will provide a powerful resource for human genetics research.

References:

- Peltonen L and McKusick VA 2001. Dissecting human disease in the postgenomic era. *Science* 291, 1224-1229
- Li 2007. Three lectures on case-control genetic association analysis. *Briefings in bioinformatics* 9: 1-13.
- Rebbeck et al 2004. Assessing the function of genetic variants in candidate gene association studies 5: 589-
- Balding D 2006. A tutorial on statistical methods for population association studies. *Nature Reviews Genetics*, 7, 781-791.

Background reading:

- Hardy et al 2009. Genomewide association studies and human disease. *NEJM* 360: 1786-.
- Kruglyak L 2008. The road to genomewide association studies. *Nature Reviews Genetics* 9: 314-
- Wang et al 2005. Genome-wide association studies: theoretical and practical concerns. *Nature Reviews Genetics* 6: 109-
- Ensenauer et al 2003. *Primer on medical genomics. Part VIII: essentials of medical genetics for the practicing physician*

In-class discussion document

- Balding D 2006. A tutorial on statistical methods for population association studies. *Nature Reviews Genetics*, 7, 781-791.

Questions: In class reading_8.pdf

Preparatory Reading:

- Laird and Lange 2006. Family-based designs in the age of large-scale gene-association studies. *Nature Reviews Genetics* 7, 385-394