

# Population Structure and Cryptic Relatedness in Genetic Association Studies

William Astle and David J. Balding<sup>1</sup>

*Abstract.* We review the problem of confounding in genetic association studies, which arises principally because of population structure and cryptic relatedness. Many treatments of the problem consider only a simple “island” model of population structure. We take a broader approach, which views population structure and cryptic relatedness as different aspects of a single confounder: the unobserved pedigree defining the (often distant) relationships among the study subjects. Kinship is therefore a central concept, and we review methods of defining and estimating kinship coefficients, both pedigree-based and marker-based. In this unified framework we review solutions to the problem of population structure, including family-based study designs, genomic control, structured association, regression control, principal components adjustment and linear mixed models. The last solution makes the most explicit use of the kinships among the study subjects, and has an established role in the analysis of animal and plant breeding studies. Recent computational developments mean that analyses of human genetic association data are beginning to benefit from its powerful tests for association, which protect against population structure and cryptic kinship, as well as intermediate levels of confounding by the pedigree.

*Key words and phrases:* Cryptic relatedness, genomic control, kinship, mixed model, complex disease genetics, ascertainment.

## 1. CONFOUNDING IN GENETIC EPIDEMIOLOGY

### 1.1 Association and Linkage

Genetic *association studies* (Clayton, 2007) are designed to identify genetic loci at which the allelic state is correlated with a phenotype of interest. The associations of interest are causal, arising at loci whose different alleles have different effects on phenotype. Even if a causal locus is not genotyped in the study, it may be possible to identify

---

William Astle is Research Associate, Centre for Biostatistics, Department of Epidemiology and Public Health, St. Mary’s Hospital Campus, Imperial College London, Norfolk Place, London, W2 1PG, UK e-mail: [wja@ic.ac.uk](mailto:wja@ic.ac.uk). David J. Balding is Professor of Statistical Genetics, Centre for Biostatistics, Department of Epidemiology and Public Health, St. Mary’s Hospital Campus, Imperial College London, Norfolk Place, London, W2 1PG, UK e-mail: [d.balding@ucl.ac.uk](mailto:d.balding@ucl.ac.uk).

<sup>1</sup>Current Address: Institute of Genetics, University College London, 5 Gower Place, London, WC1E 6BT, UK.

This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in *Statistical Science*, 2009, Vol. 24, No. 4, 451–471. This reprint differs from the original in pagination and typographic detail.

an association indirectly through a genotyped locus that is nearby on the genome. In this review we are concerned with the task of guarding against spurious associations, those which do not arise at or near a causal locus. We first introduce background material describing linkage and association studies, population structure and linkage disequilibrium, the problem of confounding by population structure and cryptic relatedness. In Section 2 we discuss definitions and estimators of the kinship coefficients that are central to our review of methods of correcting for confounding by population structure and cryptic relatedness, which is presented in Section 3. Finally, in Section 4 we present the results of a small simulation study illustrating the merits of the most important methods introduced in Section 3.

Although association designs are used to study other species, we will mainly take a human-genetics viewpoint. For example, we will focus on binary phenotypes, such as disease case/control or drug responder/nonresponder, which remain the most commonly studied type of outcome in humans, although quantitative (continuous), categorical and time-to-event traits are increasingly important. The subjects of an association study are sometimes sampled from a population without regard to phenotype, as in prospective cohort designs. However, retrospective

ascertainment of individuals on the basis of phenotype, as in case-control study designs, is more common in human genetics, and we will focus on such designs here.

*Linkage studies* (Thompson, 2007) form the other major class of study designs in genetic epidemiology. These seek loci at which there is correlation between the phenotype of interest and the pattern of transmission of DNA sequence over generations in a known pedigree. In contrast, association studies are used to search for loci at which there is a significant association between the phenotypes and genotypes of unrelated individuals. These associations arise because of correlations in transmissions of phenotypes and genotypes over many generations, but association analyses do not model these transmissions directly, whereas linkage analyses do. The relatedness of study subjects is therefore central to a linkage study, whereas the relatedness of association study subjects is typically unknown and assumed to be distant; any close relatedness is a nuisance (Figure 1).

In the last decade, association studies have become increasingly prominent in human genetics, while, although they remain important, the role of linkage studies has declined. Linkage studies can provide strong and robust evidence for genetic causation, but are limited by the difficulty of ascertaining

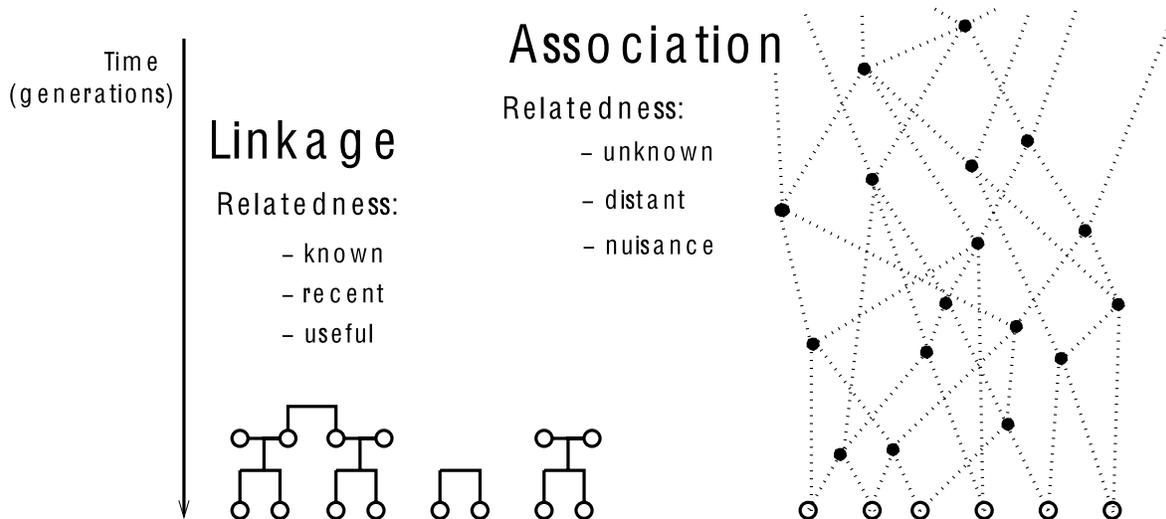


FIG. 1. Schematic illustration of differences between linkage studies, which track transmissions in known pedigrees, and population association studies which assume “unrelated” individuals. Open circles denote study subjects for whom phenotype data are available and solid lines denote observed parent-child relationships. Dotted lines indicate unobserved lines of descent, which may extend over many generations, and filled circles indicate the common ancestors at which these lineages first diverge. Unobserved ancestral lineages also connect the founders of a linkage study, but these have little impact on inferences and are ignored, whereas they form the basis of the rationale for an association analysis and constitute an important potential confounder.

enough suitable families, and by insufficient recombinations within these families to refine the location of a causal variant. When only a few hundred of genetic markers were available, lack of within-family recombinations was not a limitation. Now, cost-effective technology for genotyping  $\sim 10^6$  *single nucleotide polymorphism* (SNP) markers distributed across the genome has made possible genome-wide association studies (GWAS) which investigate most of the common genetic variation in a population, and obtain orders of magnitude finer resolution than a comparable linkage study (Morris and Cardon, 2007; Altshuler, Daly and Lander, 2008). GWAS are preferred for detecting common causal variants (say, population fraction  $> 0.05$ ), which typically have only a weak effect on phenotype, whereas linkage studies remain superior for the detection of rare variants of large effect (because these effects are more strongly concentrated within particular families).

Because genes are essentially immutable during an individual’s lifetime, and because of the independence of allelic transmissions at unlinked loci (Mendel’s Second Law), linkage studies are virtually immune to confounding. Association studies are, however, susceptible to genetic confounding, which is usually thought of as coming in two forms: *population structure* and *cryptic relatedness*. These are in fact two ends of a spectrum of the same confounder:

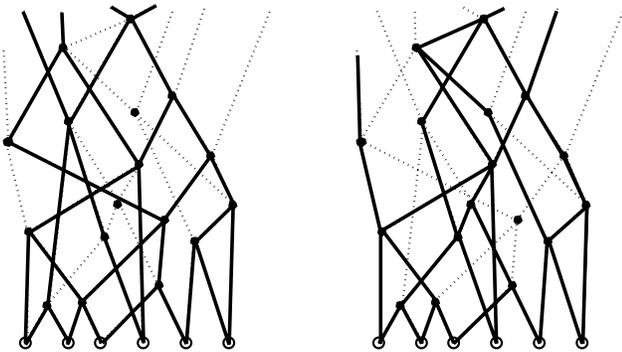


FIG. 2. Schematic illustration of the confounding role of pedigree on ancestral lineages at individual loci. Two possible single-locus lineages are shown (solid lines), each embedded in the pedigree of Figure 1 (right). Moving upwards from the study subjects (open circles), when two lineages meet at a common ancestor (filled circle), they either coalesce into a single lineage, or else they pass through different alleles of the common ancestor and do not coalesce. Dotted lines show pedigree relationships that do not contribute to the ancestry of the study subjects at this locus. Although lineages are random, they are constrained by the pedigree, features of which are therefore reflected in lineages across the genome.

the unobserved pedigree specifying the (possibly distant) relationships among the study subjects (Figure 1, right). Association studies are also susceptible to confounding if genotyping error rates vary with phenotype (Clayton et al., 2005). This can resemble a form of population structure and is not discussed further here.

We can briefly encapsulate the genetic confounding problem as follows. Association studies seek genomic loci at which differences in the genotype distributions between cases and controls indicate that their ancestries are systematically different *at that locus*. However, pedigree structure can generate a tendency for systematic ancestry differences between cases and controls at all loci not subject to strong selection. Figure 2 illustrates two possible ancestral lineages of the study subject alleles at a locus. Lineages are correlated because they are constrained to follow the underlying pedigree. For example, if the pedigree shows clustering of individuals into subpopulations, then ancestral lineages at neutral loci will tend to reflect this. The goal of correction for population structure is to allow for the confounding pedigree effects when assessing differences in ancestry between cases and controls at individual loci. In the following sections we seek to expand on this brief characterization.

## 1.2 Population Structure

Informally, a population has structure when there are large-scale systematic differences in ancestry, for example, varying levels of immigrant ancestry, or groups of individuals with more recent shared ancestors than one would expect in a *panmictic* (random-mating) population. Shared ancestry corresponds to relatedness, or *kinship*, and so population structure can be defined in terms of patterns of kinship among groups of individuals. Population structure is often closely aligned with geography, and in the absence of genetic information, stratification by geographic region may be employed to try to identify homogeneous subpopulations. However, this approach does not account for recent migration or for nongeographic patterns of kinship based on social or religious groups.

The simplest model of population structure assumes a partition of the population into “islands” (subpopulations). Mating occurs preferentially between pairs of individuals from the same island, so that the island allele fractions tend to diverge to an extent that depends on the inter-island migration

rates. An enhancement of the island model to incorporate admixture allows individual-specific proportions of ancestry arising from actual or hypothetical ancestral islands.

Below we will focus on island models of population structure, because these are simple and parsimonious models that facilitate discussion of the main ideas. Moreover, several popular statistical methods for detecting population structure and correcting association analysis for its effects have been based entirely on such models. However, human population genetic and demographic studies suggest that island models typically do not provide a good fit for human genetic data. Colonization often occurs in waves and is influenced by geographic and cultural factors. Such processes are expected to lead to clinal patterns of genetic variation rather than a partition into subpopulations (Handley et al., 2007). Modern humans are known to have evolved in Africa with the first wave of human migration from Africa estimated to have been approximately 60,000 years ago. Reflecting this history, current human genetic diversity decreases roughly linearly with distance from East Africa (Liu et al., 2006). Within Europe, Lao et al. (2008) found that the first two principal components of genome-wide genetic variation accurately reflect latitude and longitude: there is population structure at a Europe-wide level, but no natural classification of Europeans into a small number of subpopulations. Similarly, there does not appear to be a simple admixture model based on hypothetical ancestral subpopulations that can adequately capture European genetic variation, although a model based on varying levels of admixture from hypothetical “North Europe” and “South Europe” subpopulations could at least capture the latitude effect. The admixture model may be appropriate when the current population results from some intermixing following large-scale migrations over large distances, such as in Brazil or the Caribbean.

Because the term “population stratification” can imply an underlying island model, we avoid this term and adhere to “population structure,” which allows for more complex underlying demographic models.

### 1.3 Linkage Disequilibrium

In a large, panmictic population, and in the absence of selection, pairs of genetic loci that are not *tightly linked* (close together on a chromosome) are unassociated at the population level (McVean, 2007).

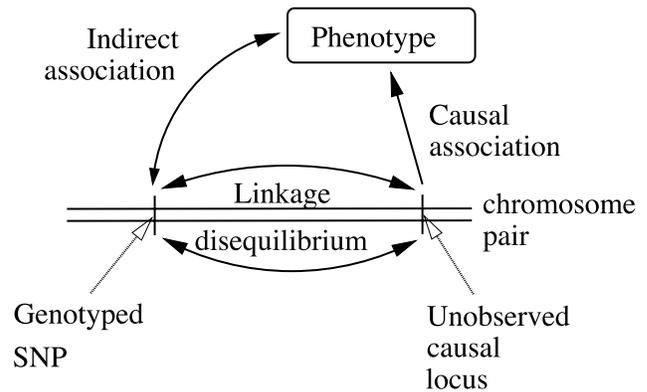


FIG. 3. Illustration of the role of linkage disequilibrium in generating phenotypic association with a noncausal genotyped marker due to a tightly-linked ungenotyped causal locus.

Such *linkage equilibrium* arises because recombination events ensure the independent assortment of alleles when they are transmitted across generations (a process sometimes called Mendelian Randomization). Conversely, because recombination is rare ( $\sim 1$  recombination per chromosome per generation), tightly linked loci are generally correlated, or in *linkage disequilibrium* (LD) in the population. This is because many individuals can inherit a linked allele pair from a remote common ancestor without an intervening recombination. Association mapping relies on LD because, even for a GWAS, only a small proportion of genetic variants are directly measured. Signals from ungenotyped causal variants can only be detected through phenotype association with a genotyped marker that is in sufficiently strong LD with the causal variant (Figure 3). LD is a double-edged sword: the stronger the LD around a causal variant, the easier it is to detect, because the greater the probability it is in high LD with at least one genotyped marker (Pritchard and Przeworski, 2001). However, in a region of high LD it is hard to fine-map a causal variant because there will be multiple highly-correlated markers each showing a similar strength of association with the phenotype.

### 1.4 Spurious Associations due to Population Structure

Unfortunately, population structure can cause LD between unlinked loci and consequently generate spurious marker-phenotype associations. For example, in the island model of population structure, if the proportion of cases among the sampled individuals varies across subpopulations, then alleles that vary in frequency across subpopulations will often show

association with phenotype. One or more such alleles may in fact be involved in phenotype determination, but standard association statistics may not distinguish them from the many genome-wide alleles with frequencies that just happen to vary across subpopulations because of differential genetic drift or natural selection. To express this another way, many alleles across the genome are likely to be somewhat informative about an individual's subpopulation of origin, and hence be predictive of any phenotype that varies across subpopulations. For example, in a large sample drawn from the population of Great Britain, many genetic variants are likely to show association with the phenotype "speaks Welsh." These will be alleles that are relatively common in Wales, which has a different population history from England (Weale et al., 2002), and do not "cause" speaking Welsh.

Under an island model, one could potentially solve the problem of spurious associations by matching for ancestry, for example, by choosing for each case a control from the same subpopulation. However, as noted above, an island model is unlikely to describe the ancestry of a human population adequately. We each have a distinct pattern of ancestry, to a large extent unknown beyond a few generations, making precise matching impractical while crude matching may be insufficient. The spouse of a case, or another relative by marriage, can provide a genetically unrelated control approximately matched for ancestry, but there are obvious limitations to this approach.

There are at least three reasons why, in an unmatched study, the phenotypes of study subjects might vary systematically with ancestry (e.g., with subpopulation in an island model). The most straightforward reason is that the disease prevalence varies across subpopulations in accordance with the frequencies of causal alleles, and the differing sample case:control ratios across subpopulations reflect the differing subpopulation prevalences. Alternatively, subpopulation prevalences may vary because of differing environmental risks. Third, *ascertainment bias* can make an important contribution to associations between ancestry and phenotype. Ascertainment bias can arise if there are differences in the sampling strategies between cases and controls that are correlated with ancestry. In the island model, this means that the sample case:control ratios across subpopulations do not reflect the subpopulation prevalences. This may happen, for example, because cases, but not controls, are sampled from clinics that over-represent particular groups.

## 1.5 Extent of the Problem

The vulnerability of association studies to confounding by population structure has been recognized for many years. In a famous example, Knowler et al. (1988) found a significant association between an immunoglobulin haplotype and type II diabetes. The study subjects were native North Americans with some European ancestry and the association disappeared after stratification by ancestry. Many commentators fail to note that Knowler et al. understood the problem and performed an appropriate analysis, so that no false association was reported: they merely noted the potential for confounding in an unstratified analysis.

Marchini et al. (2004a) concluded from a simulation study that, even in populations with relatively modest levels of structure (such as Europe or East Asia), when the sample is large enough to provide the required power, the most significant SNPs can have their  $p$ -values reduced by a factor of three because of population structure, thus exaggerating the significance of the association. Freedman et al. (2004) examined a study into prostate cancer in (admixed) African Americans and estimated a similar reduction in the smallest  $p$ -values. Another study of European-Americans found a SNP in the lactase gene significantly associated with variation in height (Campbell et al., 2005). When the subjects were stratified according to North/West or South/East European ancestry, the association disappeared. Since we expect connections among lactase tolerance, diet and height, the association could be genuine and involve different diets, but the confounding with population structure makes this difficult to establish. Helgason et al. (2005) used pedigree and marker data from the Icelandic population, and found evidence of population structure in rural areas, which would result on average in a 50% increase in the magnitude of a  $\chi^2_1$  association statistic.

Following Pritchard and Rosenberg (1999) and Gorroochurn et al. (2004), Rosenberg and Nordborg (2006) considered a general model for populations with continuous and discrete structure and presented necessary and sufficient conditions for spurious association to occur at a given locus. They defined a parameter measuring the severity of confounding under general ascertainment schemes, and showed that, broadly speaking, the case of two discrete subpopulations is worse than the cases of either more subpopulations or an admixed population. As the number of subpopulations becomes larger, the problem

of spurious association tends to diminish because the law of large numbers smoothes out correlation between disease risk and allele frequencies across subpopulations (Wang, Localio and Rebbeck, 2004).

In recent years results have been published from hundreds of GWAS into complex genetic traits (NHGRI GWAS Catalog, 2009). McCarthy et al. (2008) described the current consensus. The impact of population structure on association studies should be modest “as long as cases and controls are well matched for broad ethnic background, and measures are taken to identify and exclude individuals whose GWAS data reveal substantial differences in genetic background.” This is consistent with a report from a study of type II diabetes in UK Caucasians which estimated that population structure was responsible for only  $\sim 4\%$  inflation in  $\chi^2_1$  association statistics (Clayton et al., 2005). The Wellcome Trust Case Control Consortium (2007) study of seven common diseases using a UK population sample found fewer than 20 loci exhibiting strong geographic variation. The genome-wide distribution of test statistics suggested that any confounding effect was modest and no adjustment for population structure was made for the majority of their analyses.

In conclusion, the magnitude of the effect of structure depends on the population sampled and the sampling scheme, and well-designed studies should usually suffer only a small impact. However, most of the associated variants so far identified by GWAS have been of small effect size (NHGRI GWAS Catalog, 2009), and as study sizes increase in order to detect smaller effects, even modest structure could substantially increase the risk of false positive associations.

### 1.6 Cryptic Relatedness

Cryptic relatedness refers to the presence of close relatives in a sample of ostensibly unrelated individuals. Whereas population structure generally describes remote common ancestry of large groups of individuals, cryptic relatedness refers to recent common ancestry among smaller groups (often just pairs) of individuals. Like population structure, cryptic relatedness often arises in unmatched association studies and can have a confounding effect on inferences. Indeed, Devlin and Roeder (1999) argued that cryptic relatedness could pose a more serious confounding problem than population structure. A subsequent theoretical investigation of plausible demographic and sampling scenarios (Voight and

Pritchard, 2005) showed that the effect of cryptic relatedness in well-designed studies of outbred populations should be negligible, but it can be noticeable for small and isolated populations. Using pedigree and empirical genotype data from the Hutterite population, these authors found that cryptic relatedness reduces an association  $p$ -value of  $10^{-3}$  by a factor of approximately 4, and that the smaller the  $p$ -value the greater is the relative effect.

## 2. GENETIC RELATIONSHIPS

### 2.1 Kinship Coefficients Based on Known Pedigrees

The relatedness between two diploid individuals can be defined in terms of the probabilities that each subset of their four alleles at an arbitrary locus is *identical by descent* (IBD), which means that they descended from a common ancestral allele without an intermediate mutation. The probability that the two homologous alleles within an individual  $i$  are IBD is known as its *inbreeding coefficient*,  $f_i$ . When no genotype data are available, IBD probabilities can be evaluated from the distribution of path lengths when tracing allelic lineages back to common ancestors (Figure 2), convolved with a mutation model (Malécot, 1969). More commonly, IBD is equated with “recent” common ancestry, where “recent” may be defined in terms of a specified, observed pedigree, whose founders are assumed to be completely unrelated. In theoretical models, “recent” may be defined, for example, in terms of a specified number of generations, or since the last migration event affecting a lineage. Linkage analysis conditions on the available pedigree, and in this case the definition of IBD in terms of shared ancestry within that pedigree, and the assumption of unrelated founders, cause no difficulty. However, the strong dependence on the observed pedigree, or other definition of “recent” shared ancestry, is clearly unsatisfactory for a more general definition of relatedness.

A full description of the relatedness between two diploid individuals requires 15 IBD probabilities, one for each nonempty subset of four alleles, but if we regard the pair of alleles within each individual as unordered, then just eight identity coefficients (Jacquard, 1970) are required (Figure 4). An assumption of no within-individual IBD (no inbreeding) allows these eight coefficients to be collapsed into two (Cotterman, 1940), specifying probabilities for the two individuals to share exactly one and two alleles IBD.

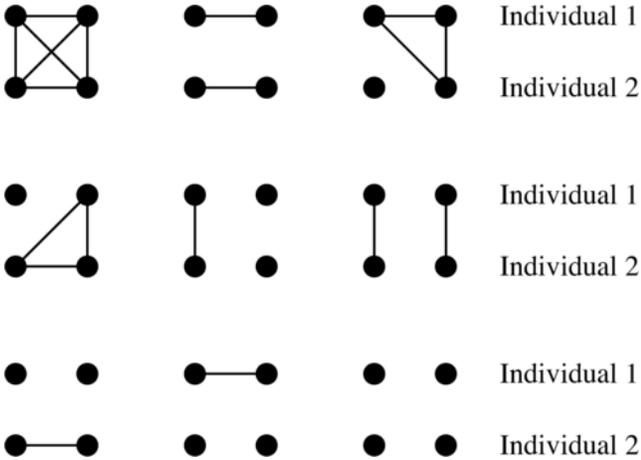


FIG. 4. Schematic illustration of the nine relatedness classes for two individuals, whose four alleles are indicated by filled circles, that are specified by the eight Jacquard identity-by-descent (IBD) coefficients. Within-individual allele pairs are regarded as unordered, and solid lines link alleles that are IBD.

Both these coefficients are required for models involving dominance, but for additive genetic models they can be reduced to a single *kinship coefficient*,  $K_{ij}$ , which is the probability that two alleles, one drawn at random from each of  $i$  and  $j$ , are IBD. Similarly,  $K_{ii}$  is the probability that two alleles, sampled with replacement from  $i$ , are IBD. Thus,  $K_{ii} = (1 + f_i)/2$ , and, in particular, the kinship of an outbred individual with itself is  $1/2$ .

The kinship matrix  $K$  of a set of individuals in a pedigree can be computed by a recursive algorithm that neglects within-pedigree mutation (Thompson, 1985).  $K$  is positive semi-definite if the submatrix of assumed founder kinships is positive semi-definite (which is satisfied if, as is typical, founders are assumed unrelated).

## 2.2 Kinship Coefficients Based on Marker Data

The advent of GWAS data means that genome-average relatedness can now be estimated accurately. It can be preferable to use these estimates in association analyses even if (unusually) pedigree-based estimates are available. There is a subtle difference between expectations computed from even a full pedigree, and realized amounts of shared genomic material. For example, if two lineages from distinct individuals meet in a common ancestor many generations in the past, then this ancestor will contribute (slightly) to the pedigree-based relatedness of the individuals but may or may not have passed any genetic material to both of them. Similarly, two pairs

of siblings in an outbred pedigree may have the same pedigree relatedness, but (slightly) different empirical relatedness (Weir, Anderson and Hepler, 2006).

Thompson (1975) proposed maximum likelihood estimates (MLEs) of the Cotterman coefficients, while Milligan (2003) made a detailed study of MLEs under the Jacquard model. These MLEs can be prone to bias when the number of markers is small and can be computationally intensive to obtain particularly from genome-wide data sets (Ritland, 1996; Milligan, 2003).

Method of moments estimators (MMEs) are typically less precise than MLEs, but are computationally efficient and can be unbiased if the ancestral allele fractions are known (Milligan, 2003). Under many population genetics models, if two alleles are not IBD, then they are regarded as random draws from some mutation operator or allele pool (Rousset, 2002), which corresponds to the notion of “unrelated.” The kinship coefficient  $K_{ij}$  is then a correlation coefficient for variables indicating whether alleles drawn from each of  $i$  and  $j$  are some given allelic type, say, A. If  $x_i$  and  $x_j$  count the numbers of A alleles (0, 1 or 2) of  $i$  and  $j$ , then

$$(2.1) \quad \text{Cov}(x_i, x_j) = 4p(1-p)K_{ij},$$

where  $p$  is the population fraction of A alleles. Thus,  $K_{ij}$  can be estimated from genome-wide covariances of allele counts. Specifically, if we write  $x$  as a column vector over individuals and let the subscript index the  $L$  loci (rather than individuals), then

$$(2.2) \quad \hat{K} = \frac{1}{L} \sum_{l=1}^L \frac{(x_l - 2p_l \mathbf{1})(x_l - 2p_l \mathbf{1})^T}{4p_l(1-p_l)}$$

is an unbiased and positive semi-definite estimator for the kinship matrix  $K$ . Entries in  $\hat{K}$  can also be interpreted in terms of excess allele sharing beyond that expected for unrelated individuals, given the allele fractions. According to Ritland (1996), who considered similar estimators and gave a generalization to loci with more than two alleles, (2.2) was first given in Li and Horvitz (1953) but only for inbreeding coefficients.

In practice, we do not know the allele fractions  $p_l$ . The natural estimators assume outbred and unrelated individuals, deviation from which can exaggerate the downward bias in the  $K_{ij}$  estimates that arises from the overfitting effect of estimating the  $p_l$  from the same data. To reduce the first problem,

one could iteratively re-estimate the  $p_l$  after making an initial estimate of  $K$  with

$$\hat{p}_l = \frac{\mathbf{1}^T \hat{K}^{-1} x_l}{\mathbf{1}^T \hat{K}^{-1} \mathbf{1}}.$$

Although the correlations arising from shared ancestry are in principle positive, because of bias arising from estimation of the  $p_l$ , off-diagonal entries of (2.2) can be negative, a property that has caused some authors to shun such estimators of  $K$  (Milligan, 2003; Yu et al., 2006; Zhao et al., 2007). Rousset (2002) also criticized the model underlying (2.1) in the context of certain population genetics models, but did not propose an alternative estimator of genetic covariance in actual populations. For our purpose, that of modeling phenotypic correlations, genotypic correlations seem intuitively appropriate and the interpretation of  $K_{ij}$  as a probability seems unimportant. Under the interpretation of  $\hat{K}_{ij}$  as excess allele sharing, negative values correspond to individuals sharing fewer alleles than expected given the allele frequencies.

Table 1 shows the probability that alleles chosen at random from each of two individuals match, that is, are *identical by state* (IBS), at a genotyped diallelic locus. The genome-wide average IBS probability can be expressed as

$$(2.3) \quad \frac{1}{2L} \sum_{l=1}^L (x_l - \mathbf{1})(x_l - \mathbf{1})^T + \frac{1}{2}.$$

If the mutation rate is low, IBS usually arises as a result of IBD, and (2.3) can be regarded as an MME of the pedigree-based kinship coefficient in the limiting case that IBS implies IBD. This estimator overcomes the problem with pedigree-based estimators of dependence on the available pedigree, but it is sensitive to recurrent mutations.

Software for computing average allele sharing (IBS) is included in popular packages for GWAS analysis such as PLINK (Purcell et al., 2007). However, because the excess allele-sharing (genotypic correlation) estimator of kinship coefficients (2.2) incorporates weighting by allele frequency, it is typically more precise than (2.3). Sharing a rare allele suggests closer kinship than sharing a common allele, because the rare allele is likely to have arisen from a more recent mutation event (Slatkin, 2002). To illustrate the increased precision of (2.2) over (2.3), we simulated 500 genetic data sets comprising 200 idealized cousin pairs (no mutation, and the alleles

not IBD from the common grandparents were independent draws from an allele pool) and 800 unrelated individuals, all genotyped at 10,000 unlinked SNPs. After rescaling to ensure the two estimators give the same difference between the mean kinship estimate of cousin pairs and mean kinship estimate of unrelated pairs, the resulting standard deviations (Table 2) are about 40% larger for the total allele sharing (IBS) estimator (2.3) than for the excess allele-sharing (genetic correlation) estimator (2.2).

The marker-based estimates of kinship coefficients discussed above do not take account of LD between markers, nor do they exploit the information about kinship inherent from the lengths of genomic regions shared between two individuals from a recent common ancestor (Browning, 2008). Hidden Markov models provide one approach to account for LD (Boehnke and Cox, 1997; Epstein, Duren and Boehnke, 2000). In outbred populations, the IBD status along a pair of chromosomes, one taken from each of a pair of individuals in a sibling, half sib or parent-child relationship, is a Markov process. However, the Markovian assumption fails for more general relationships in outbred populations. When relationships are more distant, regions of IBD will tend to cluster. For example, in the case of first cousins IBD regions will cluster into larger regions that correspond to inheritance from one of the two shared grandparents. McPeck and Sun (2000) showed how to augment the Markov model to describe the IBD process when the chromosomes correspond to an avuncular or first cousin pair. Despite the invalidity of the Markov assumption, Leutenegger et al. (2003) found that in practice it can lead to reasonable estimates for relationships more distant than first-degree.

### 3. CORRECTING ASSOCIATION ANALYSIS FOR CONFOUNDING

In this review, we seek to use kinship to illuminate connections among popular methods for protecting

TABLE 2  
*Estimated standard deviations of two kinship coefficient MMEs, after linear standardization to put the estimates on comparable scales*

Estimator	Unrelated pair	Cousin pair
Genetic correlation (2.2)	5.0	5.3
IBS (2.3)	7.3	7.2

TABLE 1

*Identity-by state (IBS) coefficients at a single diallelic locus, defined as the probability that alleles drawn at random from  $i$  and  $j$  match, which gives 0.5 in the case of a pair of heterozygotes. Another definition, based on the number of alleles in common between  $i$  and  $j$ , gives 1 for a pair of heterozygotes*

Genotype of $i$	aa	Aa	AA	aa	Aa	AA	aa	Aa	AA
Genotype of $j$	aa	aa	aa	Aa	Aa	Aa	AA	AA	AA
IBS coefficient	1	$\frac{1}{2}$	0	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$	1

association analyses from confounding. Many of these methods can be formulated within standard regression models that express the expected value of  $y_i$ , the phenotype of the  $i$ th individual, as a function of its genotype  $x_i$  at the SNP of interest:

$$(3.1) \quad g(\mathbb{E}[y_i]) = \alpha + x_i\beta,$$

where, for simplicity, we have not included covariates. Here  $g$  is a link function and  $\beta$  is a scalar or column vector of genetic effect parameters at the SNP. Often  $x_i$  counts the number of copies of a specified allele carried by  $i$ , or it can be a two-dimensional row vector that implies a general genetic model.

For a case-control study,  $g$  is typically the logit function and  $\beta$  are log odds ratios. This is a prospective model, treating case-control status as the outcome, but inferences about  $\beta$  are typically the same as for the retrospective model, which is more appropriate for case-control data (Prentice and Pyke, 1979; Seaman and Richardson, 2004). However, in some settings ascertainment effects are not correctly modeled prospectively, and it is necessary to consider retrospective models of the type

$$(3.2) \quad g(\mathbb{E}[x_i]) = \alpha + y_i\beta,$$

where  $g$  is typically the identity function.

### 3.1 Family-Based Tests of Linkage and Association (FBTLA)

The archetypal FBTLA is the transmission disequilibrium test (TDT) (Spielman, McGinnis and Ewens, 1993) for systematic differences between the genotypes of affected children and those expected under Mendelian randomization of the alleles of their unaffected parents. If an allele is directly risk-enhancing, it will be over-transmitted to cases. If not directly causal but in LD with a causal allele, it may also be over-transmitted, but in this case it must also be linked with the causal variant, since otherwise Mendelian randomization will eliminate the association between causal and tested alleles.

Thus, the TDT is a test for both association and linkage. The linkage requirement means that the test is robust to population structure, while the association requirement allows for fine-scale localization.

Parents that are homozygous at the tested SNP are uninformative and not used. Transmissions from heterozygote parents are assumed to be independent, which implies a multiplicative disease model. Let  $n_a$  and  $n_A$  denote respectively the number of a and A alleles transmitted to children by Aa heterozygote parents. If there is no linkage, each parental allele is equally likely to be transmitted, so that the null hypothesis for the TDT is

$$H_0: \mathbb{E}[n_a] = \mathbb{E}[n_A].$$

Conditional on the number of heterozygote parents  $n_a + n_A$ , the test statistic  $n_a$  has a Binomial( $n_a + n_A$ , 1/2) null distribution, but McNemar's statistic

$$(3.3) \quad \frac{(n_a - n_A)^2}{n_a + n_A},$$

which has an approximate  $\chi_1^2$  null distribution (Agresti, 2002), is widely used instead. The TDT can be derived from the score test of a logistic regression model in which transmission is the outcome variable, and the parental genotypes are predictors (Dudbridge, 2007). In Section 3.3 we outline a test which can exploit between-family as well as within-family information when it is available, while retaining protection from population structure. Tiwari et al. (2008) survey variations of the TDT in the context of a review of methods of correction for population structure.

The main disadvantages of the TDT and other FBTLA are the problem of obtaining enough families for a well powered study (particularly for adult-onset diseases) and the additional cost of genotyping: three individuals must be genotyped to obtain the equivalent of one case-control pair, and homozygous parents are uninformative. Given the availability of good analysis-based solutions to the problem of population structure (see below), the design-based solution of the FBTLA pays too high a price

for protection against spurious associations (Car- don and Palmer, 2003). However, FBTLA designs (like other linkage designs) can also be used to investigate parent-of-origin effects (Weinberg, 1999), which is not usually possible for population-based case-control studies.

### 3.2 Genomic Control

Genomic Control (GC) is an easy-to-apply and computationally fast method for reducing the inflation of test statistics caused by population structure or cryptic relatedness. It can be applied to data of any family structure or none. GC was developed (Devlin and Roeder, 1999) for the Armitage test statistic, which is asymptotically equivalent to a score statistic under logistic regression (Agresti, 2002) and, in the absence of confounding, has an asymptotic  $\chi_1^2$  null distribution. The Armitage test assumes an additive disease model, but GC has also been adapted for tests of other disease models (Zheng et al., 2005; Zheng, Freidlin and Gastwirth, 2006).

Figure 5(A) illustrates the inflation of Armitage test statistics at 2000 null SNPs simulated under an island model with admixture and ascertainment bias. This inflation could reflect many genome-wide true associations, but it is more plausible (and correct for this simulation) that the inflation is due to a combination of population structure and ascertainment bias. The figure suggests that the inflation of test statistics is approximately linear, and Devlin and Roeder argued that this holds more generally. They therefore proposed to calibrate the type I error of the Armitage test by adjusting all test statistics by a constant factor  $\lambda$ . This leaves the ranking of markers in terms of significance unchanged [Figure 5(B)], and so GC is equivalent to adjusting the significance threshold.

For most complex phenotypes, only a few genome-wide SNPs correspond to strong causal associations, with test statistics in the upper tail of the empirical distribution. Consequently, the bulk of the empirical distribution, away from the upper tail, should reflect the null distribution and can be used to estimate  $\lambda$ . Bacanu, Devlin and Roeder (2000) suggested estimating  $\lambda$  as the ratio of the empirical median to its null value ( $=0.455$ ), because the median is robust to a few large values in the upper tail. For the simulation of Figure 5, the median of the test statistics is 0.59, leading to  $\lambda = 1.31$ , a large value reflecting the strong ascertainment bias.

Setakis, Stirnadel and Balding (2006) pointed out that ascertainment bias can cause median-adjusted GC to be very conservative. Marchini et al. (2004a) had previously noticed that for strong population structure GC can be anti-conservative when the number of test statistics used to estimate  $\lambda$  is  $<100$ , and conservative when the number is  $\gg 100$ . Devlin, Bacanu and Roeder (2004) ascribed this problem to failure to account for the uncertainty in the estimate of  $\lambda$ , but Marchini et al. (2004b) noted that this may not be the most important cause of the problem (see also below). To allow for this uncertainty, Devlin, Bacanu and Roeder (2004) suggested using the mean of the test statistics to estimate  $\lambda$ , since the mean-adjusted test statistics have an  $F_{1,m}$  null distribution. In the absence of true associations, Dadd, Weale and Lewis (2009) found mean-adjusted GC to be slightly superior to median adjustment. However, the median is more robust to true positives than the mean. As a compromise, Clayton et al. (2005) proposed adjusting on a trimmed mean, discarding say the highest 5% or 10% of test statistics.

LEMMA 1. *The mean of the smallest  $100q\%$  values in a large random sample of  $\chi_1^2$  statistics has expected value*

$$\frac{1}{q}d_3(d_1^{-1}(q)),$$

where  $d_k$  is the distribution function of a  $\chi_k^2$  random variable.

PROOF. Let  $X \sim \chi_1^2$ , then

$$\begin{aligned} \mathbb{E}(X|X < d_1^{-1}(q)) &= \int_0^{d_1^{-1}(q)} x \frac{1}{q\sqrt{2\pi}} \frac{e^{-x/2}}{\sqrt{x}} dx \\ &= \int_0^{d_1^{-1}(q)} \frac{1}{q\sqrt{2\pi}} \sqrt{x} e^{-x/2} dx \\ &= \frac{1}{q}d_3(d_1^{-1}(q)). \quad \square \end{aligned}$$

A limitation of all GC methods is that they do not distinguish markers at which the pattern of association is correlated with the underlying pedigree from those at which the pedigree does not contribute to the association and so for which no adjustment should be necessary. Figure 5(B) and (C) shows that median-GC-adjustment performs similarly to PC-adjustment (see Section 3.6 below) in countering the overall inflation of test statistics, but the corrected statistics can be very different [Figure 5(D)] because

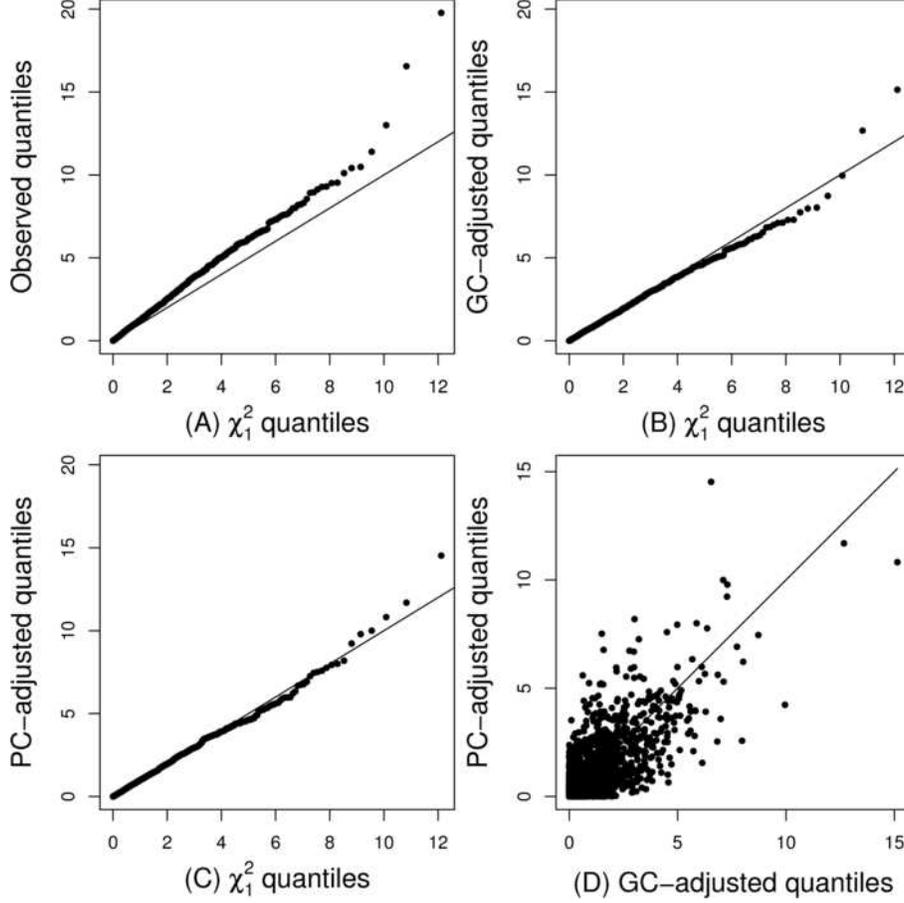


FIG. 5.  $Q$ - $Q$  plots for likelihood ratio tests of association in logistic regression (equivalent to the Armitage trend test), at 2000 null SNPs simulated under a three-island model with  $F_{ST} = 1\%$ . From Island 1 there are 200 controls and 100 cases. Each of the remaining 700 individuals is admixed, the  $i$ th individual having a proportion  $a_i$  of their ancestry from Island 2, the remainder from Island 3, where the  $a_i$  are independent and Uniform(0,1). The  $i$ th admixed individual has a probability  $0.3 + 0.5 \times a_i$  of being a case, so that case status is positively correlated with Island 2 ancestry. (A) expected versus observed quantiles, unadjusted; (B) expected versus observed after GC median-adjustment; (C) expected versus observed when the first two principal components are included as covariates; (D) GC-adjusted versus PC-adjusted quantiles.

PC-adjustment is SNP-specific. GC often shows reduced power to detect association compared to rival methods for adjusting for population structure.

In the remainder of this section we show connections between  $\lambda$  and the kinship of study subjects. The Armitage test statistic can be written as  $T^2/V$ , where  $T$  is the difference between the allele fractions in the samples of  $n_1$  cases and  $n_0$  controls,

$$T = \sum_i \left( \frac{y_i}{n_1} - \frac{1 - y_i}{n_0} \right) x_i,$$

$V$  is an estimate of the variance of  $T$ ,

$$V = \left( \frac{1}{n_0} + \frac{1}{n_1} \right) \left( \frac{1}{n} \sum_i x_i^2 - \left[ \frac{1}{n} \sum_i x_i \right]^2 \right),$$

and  $n = n_0 + n_1$ . In the following we assume retrospective ascertainment, so that the case/control status  $y$  is fixed by the study design, while the allele count  $x_i$  is random. Devlin and Roeder (1999) noticed that  $\mathbb{E}[T] = 0$ , irrespective of population structure, but that  $\text{Var}[T]$  can be inflated relative to  $V$ . In general,

$$\begin{aligned} \text{Var}[T] = \sum_{i,j} \left( \frac{y_i y_j}{n_1^2} + \frac{(1 - y_i)(1 - y_j)}{n_0^2} \right. \\ \left. - \frac{(y_i - y_j)^2}{n_1 n_0} \right) \\ \cdot \text{Cov}(x_i, x_j), \end{aligned} \quad (3.4)$$

and substituting (2.1) into (3.4) leads to

$$\text{Var}[T] = \frac{4p(1-p)}{n_0 n_1} (D + R),$$

where

$$\begin{aligned} D &= \sum_i \left( \frac{n_0}{n_1} y_i + \frac{n_1}{n_0} (1 - y_i) \right) K_{ii} \\ &\geq \min \left( \frac{n_0}{n_1}, \frac{n_1}{n_0} \right) \text{Tr}[K], \\ R &= \sum_{i \neq j} \left( \frac{n_0}{n_1} y_i y_j + \frac{n_1}{n_0} (1 - y_i)(1 - y_j) \right) K_{ij} \\ &\quad - \sum_{i \neq j} (y_i - y_j)^2 K_{ij}. \end{aligned} \tag{3.5}$$

It also follows that

$$\mathbb{E}[V] = \frac{4p(1-p)}{n_0 n_1} \left( \sum_i K_{ii} - \frac{1}{n} \sum_{i,j} K_{ij} \right), \tag{3.6}$$

which reduces to  $2p(1-p)n/n_0 n_1$  if all study subjects are outbred and unrelated. Thus, provided that

$$\frac{V}{\mathbb{E}[V]} \xrightarrow{P} 1 \quad \text{as } n \rightarrow \infty, \tag{3.7}$$

we have

$$\lambda = \mathbb{E} \left[ \frac{T^2}{V} \right] \approx \frac{\text{Var}[T]}{\mathbb{E}[V]} = \frac{D + R}{\sum_i K_{ii} - (1/n) \sum_{i,j} K_{ij}}.$$

Since  $K$  is positive semi-definite, the second summation in (3.6) is  $\geq 0$ , so that

$$\lambda \geq \frac{D + R}{\text{Tr}[K]} \geq \min \left( \frac{n_0}{n_1}, \frac{n_1}{n_0} \right) + \frac{R}{\text{Tr}[K]}.$$

The dominant quantity bounding  $\lambda$  is  $R/\text{Tr}[K]$  and since  $\text{Tr}[K] \propto n$ ,

$$\lambda \sim \frac{R}{\text{Tr}[K]} \sim R/n.$$

The large  $n$  behavior of  $\lambda$  depends on that of  $R$ . From (3.5) we see that increasing levels of kinship either among cases or among controls will tend to increase  $R$ , while greater case-control kinship tends to reduce  $R$ . In the worst-case scenario, a typical individual will be related at each degree of kinship to a fixed proportion of the study sample as  $n$  varies so that, unless the average kinship among cases and among controls is balanced by the average

case-control kinship (such as when cases and controls are matched within subpopulations under an island model),  $R \propto n^2$  and

$$\lambda \sim n,$$

which generalizes the result of Devlin and Roeder (1999). In practice, it is unclear how  $\lambda$  should scale with  $n$  in well designed GWAS studies of homogeneous populations.

The statistic  $T^2/\lambda V$  has the correct median or mean, but it will not have an asymptotic  $\chi_1^2$  null distribution unless (3.7) holds. This condition does not hold, for example, in an island model with a fixed number of islands. This fact may underlie the poor performance of GC in the settings discussed above. See Zheng et al. (2010) for a detailed discussion of variance distortion in GC.

### 3.3 Explicit Modeling of Genetic Correlations

GC is designed to correct the null distribution of a test statistic derived under a probability model which is invalid in the presence of structure. An alternative strategy is to derive a statistic from a probability model that better reflects the actual data generating process. For a retrospective study, the individual allele counts  $x_i$  should be modeled as random variables satisfying (2.1). If we are prepared to assume the higher order moments of  $x$  are small, this can be achieved with a regression model of the form (3.2) with linear link and residuals

$$x - \mathbf{1}\alpha - y\beta \sim N(0, \sigma^2 K).$$

If we assign to  $\alpha$  a diffuse Gaussian prior  $N(0, \tau^{-2})$ , with  $\tau \downarrow 0$ , and to  $\sigma^2$  the improper (Jeffreys) prior with density proportional to  $\sigma^{-2}$ , we can derive the score statistic  $T^2/V$ , with a  $\chi_1^2$  asymptotic distribution, where

$$T = y^T P x$$

and

$$(n-1)V = y^T P y \cdot x^T P x - (y^T P x)^2$$

for

$$P = K^{-1} - \frac{K^{-1} \mathbf{1} \mathbf{1}^T K^{-1}}{\mathbf{1}^T K^{-1} \mathbf{1}}. \tag{3.8}$$

If the subjects are unrelated ( $2K = I$ ), then  $T$  reduces to a comparison between the mean allele counts in cases and controls, as in the Armitage test, but the variance  $V$  is slightly smaller due to the final term. When the relatedness between study subjects

is unknown, as in a typical GWAS, the estimate  $\hat{K}$  of (2.2) may be substituted for  $K$ . A similar approach, but with a different form for  $V$ , has recently been proposed by Rakovski and Stram (2009), who point out that when the kinships are known  $T^2/V$  is equivalent to the QLS statistic of Bourgain et al. (2003).

The test described here could be used to analyze family data, either using  $K$  from the known pedigree or  $\hat{K}$  estimated from genotype data. For example, for pedigree-based  $K$  in trios of two unrelated and unaffected parents and an affected child,  $T^2$  matches (up to a constant) the numerator of (3.3).  $V$  differs slightly from the denominator of (3.3), reflecting the fact that the TDT conditions on the parental genotypes, whereas the test described here treats them as random. If the kinships are estimated from genome-wide marker data, this test can exploit the between-family as well as the within-family information, thus potentially increasing power over FBTLA, while retaining protection from population structure. Moreover, this is a very general approach, which applies for any ascertainment scheme and degree of relatedness or population structure among study subjects. Thus, if a researcher were unaware of the TDT, but applied the retrospective regression model to family trio data, s/he would automatically “invent” a test similar to the TDT but with potentially superior properties.

### 3.4 Structured Association

Structured association (SA) methods are based on the island model of population structure, and assume that the ancestry of each individual is drawn from one or more of the “islands.” Popular software packages include ADMIXMAP (Hoggart et al., 2003) and STRUCTURE/STRAT (Pritchard and Donnelly, 2001; Falush, Stephens and Pritchard, 2003). These approaches model variation in ancestral subpopulation along a chromosome as a Markov process. Stratified tests for association (Clayton, 2007), such as the Mantel–Haenszel test, can then be performed to combine signals of association across subpopulations. More generally, a logistic regression model of the form (3.1) can be employed, with admixture proportions (one for each subpopulation) entering as covariates.

Similar to GC, SA methods can be effective using only  $\sim 10^2$  SNPs, but unlike GC, they can be computationally intensive, although a simplified and

fast version of SA is implemented in PLINK (Purcell et al., 2007). The number of subpopulations can be estimated from the data by optimizing a measure of model goodness of fit, but this increases the computational burden and there is usually no satisfactory estimate because, as we noted above in Section 1.2, the island model is not well suited to most human populations. Indeed, ADMIXMAP was primarily designed for admixture mapping, in which the genomes of admixed individuals are scanned for loci at which cases show an excess of ancestry from one of the founder populations (McKeigue, 2007). Because of the limited number of generations since the admixture event, this approach has features in common with linkage as well as association study designs.

### 3.5 Regression Control

Wang, Localio and Rebbeck (2005) showed that it is possible to control for population structure within a logistic regression model of the form (3.1) by including among the covariates the genotype at a single marker that is informative about ancestry. Setakis, Stirnadel and Balding (2006) proposed using a set of  $\sim 10^2$  widely-spaced SNPs, which are assumed to be noncausal (in practice, a “random” set of SNPs). These null SNPs are informative about the underlying pedigree, which we have argued forms the basis of the problem of inflation of test statistics due to population structure. Including these SNPs as regression covariates while testing a SNP of interest should eliminate most or all of the pedigree (population structure) effect.

Setakis, Stirnadel and Balding (2006) suggest two standard procedures to avoid overfitting the SNP covariates: a backward (stepwise) selection and a shrinkage penalty approach. In the absence of ascertainment bias, both methods performed similarly to GC and SA, while being computationally fast and allowing the flexibility of the regression framework. With ascertainment bias, the regression control approach substantially outperformed GC.

Another approach (Epstein, Allen and Satten, 2007), which is also related to propensity score methods, uses ancestry-informative SNPs to create a risk score, stratifies study subjects according to this score, and performs a stratified test of association.

### 3.6 Principal Component Adjustment

Zhang, Zhu and Zhao (2003) proposed controlling for population structure in quantitative trait association analysis by including principal components

(PCs) of genome-wide SNP genotypes as regression covariates. Price et al. (2006) presented a similar method, focusing on its application to case-control GWAS. PC regression is similar to the regression control approach of Setakis, Stirnadel and Balding (2006), but minimizes overfitting by using only a few linear combinations of SNPs (the PCs), rather than a larger number of individual SNPs. However, many more SNPs (typically  $\sim 10^4$ ) are required in the PC-based approach.

Let  $X$  denote a matrix with  $n$  rows corresponding to individuals and  $L$  columns corresponding to SNPs. Genotypes are initially coded as allele counts (0, 1 or 2) but are then standardized to have zero mean and unit variance. Then the  $n \times n$  matrix  $XX^T/L$  is the estimated kinship matrix  $\hat{K}$  introduced at (2.2). Since  $\hat{K}$  is symmetric and positive semi-definite, it has an eigenvalue decomposition

$$(3.9) \quad \hat{K} = \frac{1}{L}XX^T = v\Lambda v^T,$$

where the columns of  $v$  are the eigenvectors, or PCs, of  $\hat{K}$ , while  $\Lambda$  is a diagonal matrix of corresponding (nonnegative) eigenvalues in decreasing order.

Standard principal components analysis uses the  $L \times L$  matrix  $X^T X$  specifying the empirical correlations between the columns of the design matrix. Here, the variables of interest are the individuals, corresponding to rows, and, hence, we focus on  $\hat{K} = XX^T/L$ . However, because  $X$  is column-standardized, and not row-standardized,  $\hat{K}$  is not an empirical correlation matrix. In particular, the diagonal entries of  $X^T X$  are all one, whereas the diagonal entries of  $\hat{K}$  vary over individuals according to their estimated inbreeding coefficient (Section 2).

To maximize the empirical variance of  $X^T v_1$ , the first PC  $v_1$  will typically be correlated with many SNPs. For example, in a two-island model, it will have greatest correlation with the SNPs whose allele fractions are most discrepant between the two islands. Thus,  $v_1$  acts as a strong predictor of island membership, and can also identify admixed individuals (intermediate scores). More generally, in an  $S$ -island model the first  $S - 1$  PCs predict island memberships and individual admixture proportions (Patterson, Price and Reich, 2006) (see Figure 6 for an illustration with  $S = 3$ ). The subsequent PCs represent the within-island pedigree effects. Cryptic kinship typically generates weaker LD than large-scale population structure, and the effects of small groups of close (even first-degree) relatives are usually not reflected in the leading PCs.

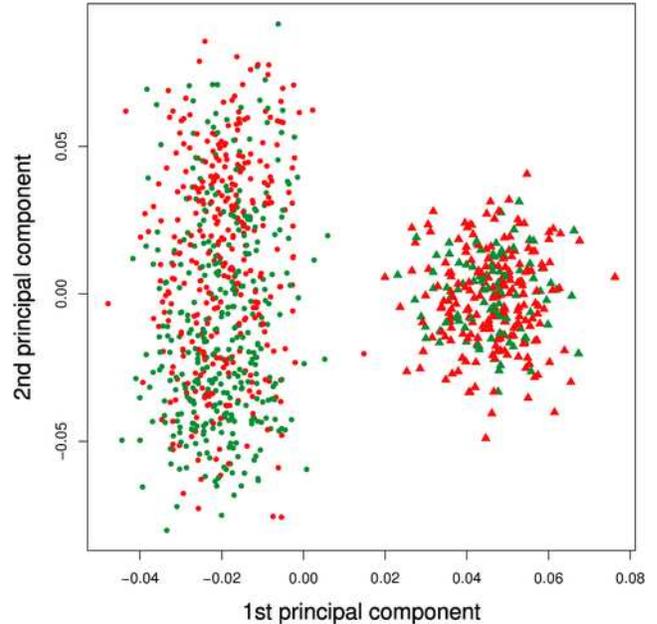


FIG. 6. First two principal component scores for the cases (green) and controls (red) of the simulation underlying Figure 5. Triangles indicate individuals from Island 1, and circles the admixed individuals with ancestry from Islands 2 and 3.

Tightly linked SNPs tend to be in high LD with each other, and sometimes one or more of the leading PCs will be dominated by large LD blocks. Since such blocks are genomically local, they convey little if any information about population structure. One way to avoid this problem is to filter GWAS SNPs prior to extracting the PCs, to exclude one in each pair of high LD SNPs.

As for structured association and regression control, the idea motivating PC adjustment is that if a correlation between the phenotype and the tested SNP can be partly explained by measures of ancestry, here PCs, then including these as regression covariates prevents that part of the signal from contributing to the test statistic. In particular, in the case of linear regression, only the components of the phenotype and genotype vectors that are orthogonal to the PCs included in the model contribute to the test statistic. This should protect against spurious associations, provided that sufficient PCs are included in the model to explain potentially confounding structure. For example, in the population of Figure 6,  $v_1$  and  $v_2$  can, with accuracy, jointly predict the proportion of an individual's ancestry arising from each of the three subpopulations, and because of the varying case-control ratios across the subpopulations, they can also to some extent predict

case-control status. Including the PCs as covariates in the regression model discards information, for example, indicating that alleles common in the high-risk subpopulation are more likely to be causal. Because of the danger of ascertainment bias in retrospective studies, such information may be dangerous and it is safer to discard it, but if ascertainment is not a problem, for example, in a prospective study, discarding this information is inefficient.

For computational reasons the EIGENSTRAT software (Price et al., 2006), which implements PC adjustment, does not include PCs as logistic regression covariates, but instead uses a linear adjustment of both phenotypes and genotypes. Such an adjustment is valid only under the assumption that the  $y_i$  form a homoskedastic sample (Agresti, 2002, page 120) and should be reasonable if the sample case:control ratio is not too far from 1, and the effect sizes are small.

By default, the EIGENSTRAT software includes the first ten PCs. Patterson, Price and Reich (2006) proposed a test to determine whether the lead eigenvalue of  $\hat{K}$  is significantly larger than one would expect under a null model, but it remains unclear what significance threshold for this test might be appropriate, if any, for protecting association test statistics from inflation. As for any other regression covariate, there is an argument for only including a PC in the model if it shows an association with the phenotype (Novembre and Stephens, 2008; Lee, Wright and Zou, 2010). Experience seems to suggest that between 2 and 15 PCs are typically sufficient, and in large studies for which  $n$  may be several thousand, these will correspond to a small loss of total genotypic information.

While the intuition motivating PC-adjustment is valid under an island model, protection from population structure effects is not guaranteed under more complicated and realistic models of population structure. In particular, inflation of test statistics due to cryptic relatedness is not ameliorated by PC adjustment. Moreover, if leading PCs reflect genome-local effects, PC adjustment could lose valuable information and lead to true effects being missed.

### 3.7 Mixed Regression Models

We assume here a quantitative phenotype with  $g$  the identity link. The linear mixed model (MM) extends (3.1) by including for each individual  $i$  a latent variable  $\delta_i$  such that

$$(3.10) \quad \mathbb{E}[y_i | \delta_i] = \alpha + x_i \beta + \delta_i.$$

The value of  $\delta_i$  is interpreted as a polygenic contribution to the phenotype, due to many small, additive, genetic effects distributed across the genome. In animal breeding genetics, the equivalent term is referred to as the *breeding value*. The additive assumption seems to be well supported for traits with a complex genetic basis (Hill, Goddard and Visscher, 2008), although it is also possible to include latent variables corresponding to dominance effects. Under the additive polygenic assumption, the variance-covariance structure of  $\delta$  is proportional to the correlation structure of genotypes coded as allele counts, which from (2.1) is proportional to  $K$ , the kinship matrix. Hence, we assume

$$(3.11) \quad \delta \sim N(0, 2\sigma^2 h^2 K),$$

where  $h^2 \in [0, 1]$  is the *narrow sense heritability*, and is defined as the proportion of phenotypic variation that can be attributed to additive polygenic effects. The residuals are assumed to satisfy

$$y_i - \alpha - x_i \beta - \delta_i \sim N(0, \sigma^2(1 - h^2)I).$$

The origins of this linear MM lie in the partitioning by Fisher (1918) of the variance of a quantitative trait into independent genetic and environmental components, and derivation of the genetic correlation of trait values of a pair of relatives assuming Mendelian inheritance (Gianola, 2007). It is conventional to introduce the 2 in (3.11) because  $2K$  reduces to  $I$  in the limiting case of completely unrelated and completely outbred individuals, in which case  $h^2$  becomes inestimable.

The model (3.10) has long been used for mapping quantitative trait loci in outbred pedigrees, using a pedigree-based kinship matrix (Höschle, 2007). Yu et al. (2006), who were interested in association mapping in maize, were first to suggest using the same model to correct association analysis for population structure, but with  $K$  estimated from marker data. In fact, Yu et al. also include in their model additional population structure terms, namely, the ancestral proportions estimated by the STRUCTURE software (Section 3.4). Zhao et al. (2007) used the same model to analyze an *Arabidopsis thaliana* data set and, like Yu et al., found that the additional terms improved the fit of the model. However, neither set of authors formally assess the improvements in fit which, by visual inspection of the genome-wide  $p$ -value distributions, seem modest. In principle,  $K$  already includes population structure information, making the additional terms redundant. However,

the structure terms provide a low-dimensional summary of key features of  $K$  which are likely to be better estimated than individual kinships, and this may generate some advantage to including the structure terms in the regression model as well as  $K$ . Typically,  $\sim 10^5$  SNPs are required for adequate estimation of  $K$  in human populations, more than are required for estimation of PCs, but this number is usually available from a GWAS.

Kang et al. (2008) have developed the software package EMMA for fast inference in linear MMs using a likelihood ratio test. Alternatively, for very large samples one can use the score test which is computationally faster because it only requires parameter estimates under the null hypothesis ( $\beta = 0$ ). Another fast method for inference in mixed models, GRAMMAR, has been proposed by Aulchenko, de Koning and Haley (2007). Although GRAMMAR is faster than EMMA, it is an approximate method and the authors found that it could be conservative and hence have reduced power. GRAMMAR uses the mixed model to predict the phenotype under the null hypothesis,

$$\hat{y} = \hat{\alpha} + \hat{\delta},$$

where  $\hat{\delta}$  is the best linear unbiased predictor (BLUP) of  $\delta$ , which is equivalent to the empirical Bayes estimate for  $\delta$  with prior (3.11) (Robinson, 1991). This prediction only needs to be made once for the whole data set. The next step is to use the residuals from the prediction as the outcome in a linear regression,

$$(3.12) \quad y - \hat{y} = \mathbf{1}\mu + x\beta + \varepsilon,$$

and test the parameter  $\beta$  for each SNP. An assumption underlying (3.12) is that the residuals are independent and identically distributed, which is strictly false. Indeed, both  $\mathbb{E}(y - \hat{y})$  and  $\text{Var}(y - \hat{y})$  are functions of  $K$  unless  $h^2 = 0$ , which may be the reason that an additional GC-style variance inflation correction is often required to calibrate the GRAMMAR type I error rate.

Note that (3.10) can be reparametrized as

$$\mathbb{E}[y_i|\gamma] = \alpha + x_i\beta + v_i\gamma,$$

with

$$\gamma \sim N(0, 2\sigma^2 h^2 \Lambda),$$

where  $\Lambda$  and  $v_i$  are defined above at (3.9). Thus, the MM approach uses the same latent variables as PC adjustment but deals differently with the vector  $\gamma$  of nuisance parameters. From a Bayesian point of

view, both methods put independent priors on the components of  $\gamma$ . Using  $k$  PCs as regression covariates can be viewed as assigning to each of the  $n - k$  trailing components of  $\gamma$  a prior with unit mass at zero, while each of the  $k$  leading components receives a diffuse prior. These assignments imply certainty that the polygenic component of the phenotype is fully captured by the first  $k$  PCs. In contrast, the MM approach puts a Gaussian prior on each component of  $\gamma$ , with variances proportional to the corresponding eigenvalues.

## 4. SIMULATIONS

### 4.1 Case-Control Studies Without Ascertainment Bias

We show here the results of a small simulation study designed to illustrate the merits of some of the methods introduced above for correcting GWAS analysis for population structure. Although we have criticized the island model as unrealistic, it remains the most convenient starting point and we use it below, the only complication considered here being ascertainment bias. More extensive and realistic simulations will be published elsewhere.

We simulated data from 500 case-control studies, each with 1000 cases and 1000 controls drawn from a population of 6000 individuals partitioned into three equal-size subpopulations. Ancestral minor allele fractions were Uniform[0.05, 0.5] for all 10,000 unlinked SNPs. For each SNP, we drew subpopulation allele fractions from the beta-binomial model described in Balding and Nichols (1995). Under this model, a marker with ancestral population allele fraction  $p$  has subpopulation allele fractions that are independent draws from

$$\text{Beta}\left(\frac{1-F}{F}p, \frac{1-F}{F}(1-p)\right),$$

where  $F$  is Wright's  $F_{ST}$ , a measure of population divergence (Balding, 2003). In order to discriminate among the methods, we simulated a high level of population structure,  $F = 0.1$ , which is close to between-continent levels of human differentiation; this is larger than is typical for a well-designed GWAS, but some meta-analyses may include populations at this level of differentiation. The studies simulated here are relatively small by the standards of current GWAS, and larger studies will be affected by less pronounced structure than that simulated here.

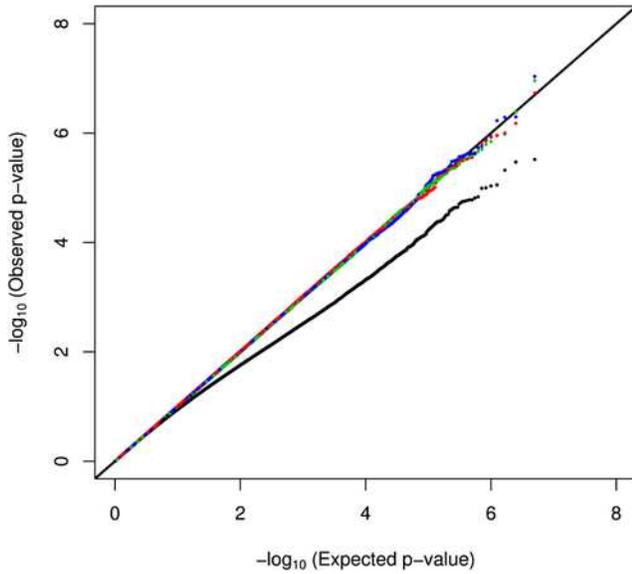


FIG. 7. *Expected versus observed  $-\log_{10}$  (p-values) for four test statistics, GC (black), PC10 (blue), MM (red) and MCP (green) evaluated at 5 million SNPs (10,000 per dataset) simulated under the null of no association in population case-control studies (see text for details of the simulation).*

We simulated the disease phenotype under a logistic regression model, with 20 SNP markers each assigned allelic odds ratio 1.18. The population disease prevalence was 0.18. We performed tests for disease-phenotype association for all the markers in each data set, using median-adjusted GC, principal component adjustment with 10 PCs (PC10) and the likelihood ratio test from the linear mixed model (MM). Note that the PC10 and MM approaches apply a linear regression model to binary outcome data, which (as noted in Section 3.6) should be reasonable if the case:control ratio is not extreme and the effect sizes are small. We also performed the score test described in Section 3.3. This retrospective model is consistent with the case-control ascertainment, although in fact the resulting score statistic is symmetric in  $x$  and  $y$ . We call this test MCP to stand for Multivariate Gaussian model Conditional on Phenotype.

The PC10, MM and MCP approaches all require specification of  $K$ , the kinship matrix of the study subjects. Figure 7 shows the observed and expected  $p$ -values at the null markers, aggregated over the 500 simulations, using  $\hat{K}$  estimated via genetic correlations (2.2). We see that the type I error is well calibrated for all the methods except GC, which has a conservative null distribution. We repeated the analysis using the true  $K$  used for the simulation and the results are similar (not shown). GC is conservative because in this extreme scenario with large  $F_{ST}$

and a small number of subpopulations, the assumption that the test statistic asymptotically follows a linear-inflated  $\chi_1^2$  distribution fails. The condition (3.7) is not satisfied in the beta-binomial model with fixed  $F_{ST}$  unless the number of subpopulations increases with  $n$ .

In order to compare power across the methods, we plotted ROC curves for the statistics from the four methods using both the true  $K$  and the estimated  $\hat{K}$  (Figure 8). GC has lower power than the other three methods, even though the ROC calibrates its bad false positive rate. When  $K$  is used, the MCP and MM approaches are equally powerful, and more powerful than the other two methods, because both exploit the between-subpopulation information. When  $\hat{K}$  is used both the MM and MCP methods lose their power advantage over PC10, which may be due to the sampling error in the eigenvectors of  $\hat{K}$ . PC correction uses only the leading eigenvectors to adjust the analysis; these are less affected by noise and for an island model they contain all the population structure signal. For an actual GWAS, the MM and MCP methods should show performance somewhere intermediate between the cases considered here, because many more than 10,000 SNPs are available to estimate  $K$ .

#### 4.2 Case-Control Studies With Ascertainment Bias

We repeated the simulation studies described above but with ascertainment bias. Specifically, each simulated study sampled 50 controls from two of the three subpopulations and 900 controls from the remaining subpopulation. This corresponds to a scenario in which investigators are forced to search widely across subpopulations to obtain a sufficient number of cases for their study, but are able to recruit most controls from the local subpopulation. The case:control ratio varies dramatically over subpopulations in this scenario, so that subpopulation allele frequencies are strong predictors of case-control status.

Once again, PC10, MM and MCP all have good control of type I error, while GC is dramatically conservative (not shown). We again plotted ROC curves for the statistics from the four methods using both  $K$  and  $\hat{K}$  (Figure 9). GC shows almost no power in these simulations. When  $K$  is used the MCP and MM approaches are equally powerful and more powerful than PC10. Their power advantage over PC10 is more substantial than in the previous

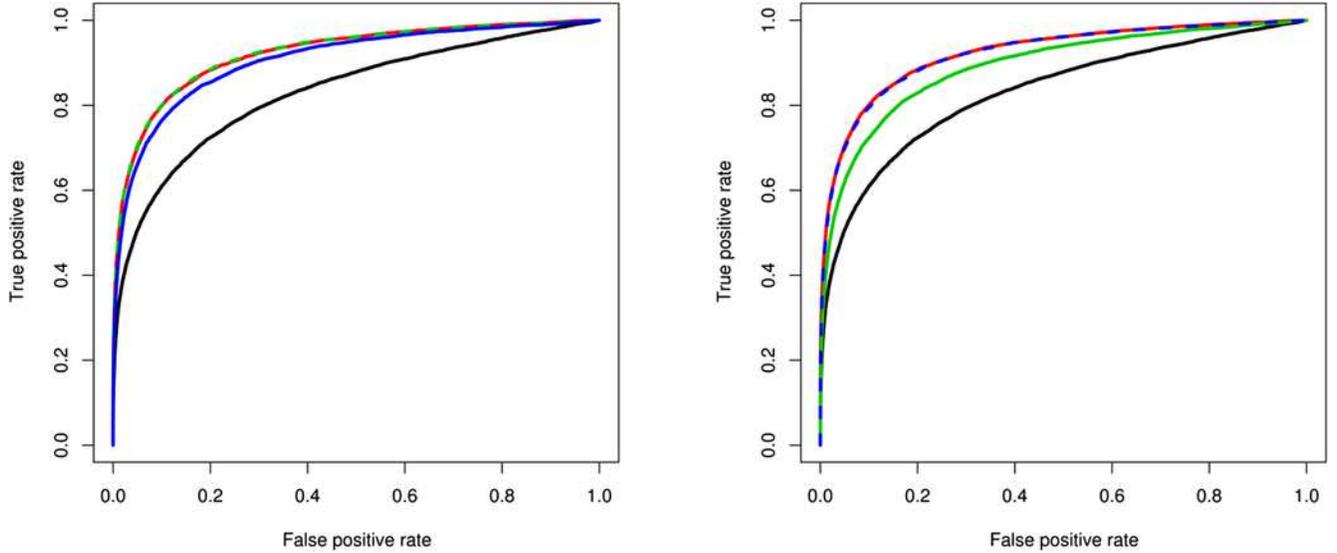


FIG. 8. ROC plots comparing true and false-positive rates of GC (black), PC10 (blue), MM (red) and MCP (green) estimated from data aggregated over 500 simulated retrospective population case-control studies. Left: true kinships,  $K$ ; right: estimated kinships,  $\hat{K}$ .

scenario without ascertainment bias, because here the leading eigenvectors of  $K$  are stronger predictors of case-control status. When  $\hat{K}$  is used the MCP and MM methods again lose some power compared with PC10 and again the MCP method suffers a greater loss than the MM test.

## 5. SUMMARY

The theme of our review has been the unifying role of the matrix of kinship coefficients,  $K$ , and its estimate  $\hat{K}$  defined at (2.2). We view population structure and cryptic kinship as the extremes of the same confounder, the latent pedigree, and  $\hat{K}$  as a good summary of the pedigree for use in adjusting association analyses. We have also argued that, whereas methods are often tested under an island model of population structure, these models do not provide realistic descriptions of relatedness in human populations.

The median-adjusted Genomic Control (GC) is simple to apply and, for association studies with moderate sample sizes and small amounts of within sample relatedness, it is a satisfactory method for protecting against confounding, which requires relatively few SNPs ( $\sim 10^2$ ). When the study sample is drawn from a population with a few, distinct subpopulations, GC should be used with caution because the  $\chi^2$  approximation may fail. Our simulations also confirmed previous reports that GC is very

sensitive to ascertainment bias. Structured association and regression adjustment may also be used with relatively few SNPs, and the latter has important advantages over GC.

Linear mixed models (MM) and the multivariate Gaussian model conditioning on phenotype (MCP) explicitly model genetic correlations using  $K$ , and are respectively appropriate for prospective and retrospective studies with genome-wide SNP data. They are particularly suited to modeling complex patterns of kinship, including intermediate scenarios between close relationships and large scale structure, which can arise in plant genetics, human population isolates and in animal breeding studies. Previous computational limitations have largely been overcome in recent years. These models also provide powerful and computationally efficient methods for analyzing family data of any structure, and combinations of families and apparently unrelated individuals.

Principal components (PC) adjustment in effect eliminates from test statistics the part of the phenotypic signal that can be predicted from large scale population structure. In particular, if natural selection leads to large allele frequency differences across subpopulations at a particular SNP, and case-control ratios vary across subpopulations, then spurious associations can arise that PC adjustment will control, because the SNP genotypes are strongly correlated with subpopulation, whereas the MM and MCP methods will not. On the other hand, the MM

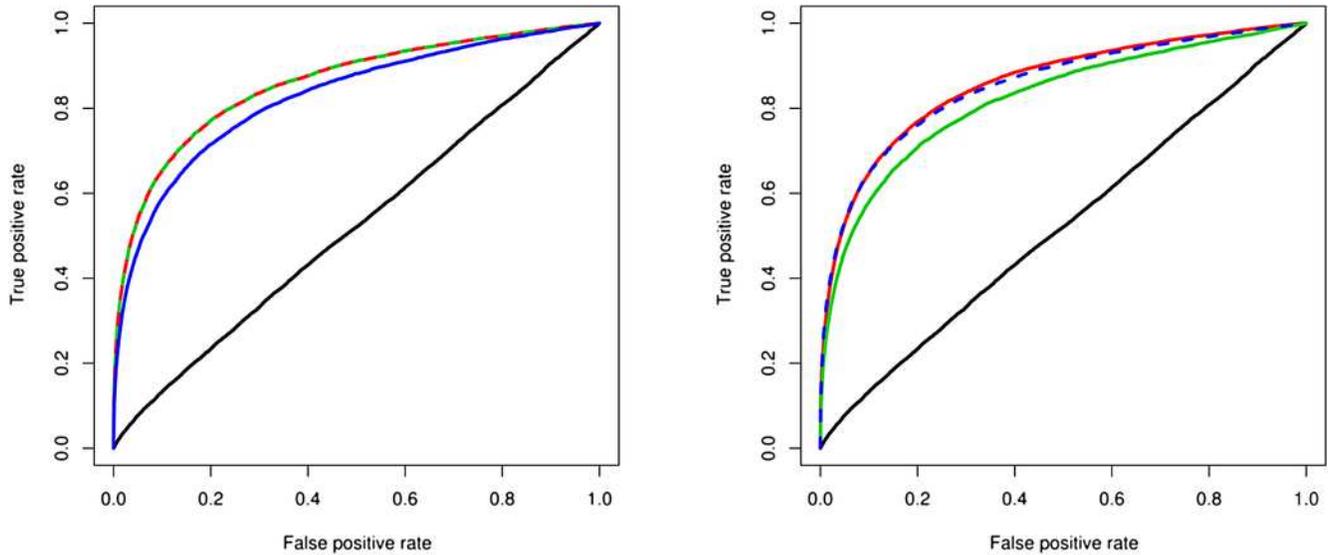


FIG. 9. ROC plots comparing true and false-positive rates of GC (black), PC10 (blue), MM (red) and MCP (green) estimated from data aggregated over 500 simulated population case-control studies with biased ascertainment of controls. Left: true kinships,  $K$ ; right: estimated kinships,  $\hat{K}$ .

and MCP methods can gain power over PC adjustment because they explicitly model the phenotype-genotype correlations induced by relatedness and genetic drift. For example, they should provide better power than PC adjustment when analyzing data from human population isolates which are homogeneous for environmental exposures.

There are probably better ways to extract the signals of population structure from an estimate of  $K$  than those considered here. Selection of the first few principal components of  $\hat{K}$  can be viewed as a form of signal denoising, but PC regression adjustment does not in general optimize power. It may be possible to adapt the MM approach to maintain the power advantage over PC while reducing noise, for example, by smoothing or truncating the lower order eigenvalues of  $\hat{K}$ .

## ACKNOWLEDGMENTS

We gratefully acknowledge helpful discussions from many colleagues, including Yurii Aulchenko, David Clayton, Clive Hoggart, Chris Holmes, Matti Pirinen, Sylvia Richardson and Jon White, and, in particular, Francois Balloux, Dan Stram and Mike Weale for helpful comments on an early draft. This work was supported in part by the UK Medical Research Council and GlaxoSmithKline.

## REFERENCES

- AGRESTI, A. (2002). *Categorical Data Analysis*, 2nd ed. Wiley, New York. [MR1914507](#)
- ALTSHULER, D., DALY, M. J. and LANDER, E. S. (2008). Genetic mapping in human disease. *Science* **322** 881–888.
- AULCHENKO, Y. S., DE KONING, D.-J. and HALEY, C. (2007). Genomewide rapid association using mixed model and regression: A fast and simple method for genomewide pedigree-based quantitative trait loci association analysis. *Genetics* **177** 577–585.
- BACANU, S. A., DEVLIN, B. and ROEDER, K. (2000). The power of genomic control. *Am. J. Hum. Genet.* **66** 1933–1944.
- BALDING, D. J. (2003). Likelihood-based inference for genetic correlation coefficients. *Theor. Popul. Biol.* **63** 221–230.
- BALDING, D. J. and NICHOLS, R. A. (1995). A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity. *Genetica* **96** 3–12.
- BOEHNKE, M. and COX, N. J. (1997). Accurate inference of relationships in sib-pair linkage studies. *Am. J. Hum. Genet.* **61** 423–429.
- BOURGAIN, C., HOFFJAN, S., NICOLAE, R., NEWMAN, D., STEINER, L., WALKER, K., REYNOLDS, R., OBER, C. and MCPEEK, M. S. (2003). Novel case-control test in a founder population identifies P-selectin as an atopy-susceptibility locus. *Am. J. Hum. Genet.* **73** 612–626.
- BROWNING, S. R. (2008). Estimation of pairwise identity by descent from dense genetic marker data in a population sample of haplotypes. *Genetics* **178** 2123–2132.
- CAMPBELL, C. D., OGBURN, E. L., LUNETTA, K. L., LYON, H. N., FREEDMAN, M. L., GROOP, L. C., ALTSHULER, D., ARDLIE, K. G. and HIRSCHHORN, J. N. (2005). Demonstrating stratification in a European American population. *Nat. Genet.* **37** 868–872.

- CARDON, L. R. and PALMER, L. J. (2003). Population stratification and spurious allelic association. *Lancet* **361** 598–604.
- CLAYTON, D. (2007). Population association. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 1264–1237. Wiley, Chichester. [MR2391785](#)
- CLAYTON, D. G., WALKER, N. M., SMYTH, D. J., PASK, R., COOPER, J. D., MAIER, L. M., SMINK, L. J., LAM, A. C., OVERTON, N. R., STEVENS, H. E., NUTLAND, S., HOWSON, J. M. M., FAHAM, M., MOORHEAD, M., JONES, H. B., FALKOWSKI, M., HARDENBOL, P., WILLIS, T. D. and TODD, J. A. (2005). Population structure, differential bias and genomic control in a large-scale, case-control association study. *Nat. Genet.* **37** 1243–1246.
- COTTERMAN, C. (1940). A calculus for statistico-genetics. Dissertation, Ohio State Univ.
- DADD, T., WEALE, M. E. and LEWIS, C. M. (2009). A critical evaluation of genomic control methods for genetic association studies. *Genet. Epidemiol.* **33** 290–298.
- DEVLIN, B., BACANU, S.-A. and ROEDER, K. (2004). Genomic control to the extreme. *Nat. Genet.* **36** 1129–1130; author reply 1131.
- DEVLIN, B. and ROEDER, K. (1999). Genomic control for association studies. *Biometrics* **55** 997–1004.
- DUDBRIDGE, F. (2007). Family-based association. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 1264–1285. Wiley, Chichester. [MR2391785](#)
- EPSTEIN, M. P., ALLEN, A. S. and SATTEN, G. A. (2007). A simple and improved correction for population stratification in case-control studies. *Am. J. Hum. Genet.* **80** 921–930.
- EPSTEIN, M. P., DUREN, W. L. and BOEHNKE, M. (2000). Improved inference of relationship for pairs of individuals. *Am. J. Hum. Genet.* **67** 1219–1231.
- FALUSH, D., STEPHENS, M. and PRITCHARD, J. K. (2003). Inference of population structure using multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics* **164** 1567–1587.
- FISHER, R. (1918). The correlation between relatives on the supposition of Mendelian inheritance. *Transactions of the Royal Society of Edinburgh* **52** 399–433.
- FREEDMAN, M. L., REICH, D., PENNEY, K. L., McDONALD, G. J., MIGNAULT, A. A., PATTERSON, N., GABRIEL, S. B., TOPOL, E. J., SMOLLER, J. W., PATO, C. N., PATO, M. T., PETRYSHEN, T. L., KOLONEL, L. N., LANDER, E. S., SKLAR, P., HENDERSON, B., HIRSCHHORN, J. N. and ALTSHULER, D. (2004). Assessing the impact of population stratification on genetic association studies. *Nat. Genet.* **36** 388–393.
- GIANOLA, D. (2007). Inferences from mixed models in quantitative genetics. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) 678–717. Wiley, Chichester. [MR2391785](#)
- GORROOCHURN, P., HODGE, S. E., HEIMAN, G. and GREENBERG, D. A. (2004). Effect of population stratification on case-control association studies. ii. False-positive rates and their limiting behavior as number of subpopulations increases. *Hum. Hered.* **58** 40–48.
- HANDLEY, L. J. L., MANICA, A., GOUDET, J. and BALLOUX, F. (2007). Going the distance: Human population genetics in a clinal world. *Trends Genet.* **23** 432–439.
- HELGASON, A., YNGVADÓTTIR, B., HRAFNKELSSON, B., GULCHER, J. and STEFÁNSSON, K. (2005). An icelandic example of the impact of population structure on association studies. *Nat. Genet.* **37** 90–95.
- HILL, W. G., GODDARD, M. E. and VISSCHER, P. M. (2008). Data and theory point to mainly additive genetic variance for complex traits. *PLoS Genet.* **4** e1000008.
- HOGGART, C. J., PARRA, E. J., SHRIVER, M. D., BONILLA, C., KITTLES, R. A., CLAYTON, D. G. and MCKEIGUE, P. M. (2003). Control of confounding of genetic associations in stratified populations. *Am. J. Hum. Genet.* **72** 1492–1504.
- HÖSCHELE, I. (2007). Mapping quantitative trait loci in outbred pedigrees. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **1** 678–717. Wiley, Chichester.
- JACQUARD, A. (1970). *Structures Génétiques des Populations*. Masson & Cie, Paris. [MR0421723](#)
- KANG, H. M., ZAITLEN, N. A., WADE, C. M., KIRBY, A., HECKERMAN, D., DALY, M. J. and ESKIN, E. (2008). Efficient control of population structure in model organism association mapping. *Genetics* **178** 1709–1723.
- KNOWLER, W. C., WILLIAMS, R. C., PETTITT, D. J. and STEINBERG, A. G. (1988). Gm3;5,13,14 and type 2 diabetes mellitus: An association in American Indians with genetic admixture. *Am. J. Hum. Genet.* **43** 520–526.
- LAO, O., LU, T. T., NOTHNAGEL, M., JUNGE, O., FREITAG-WOLF, S., CALIEBE, A., BALASCAKOVA, M., BERTRAN-PETIT, J., BINDOFF, L. A., COMAS, D., HOLMLUND, G., KOUVATSI, A., MACEK, M., MOLLET, I., PARSON, W., PALO, J., PLOSKI, R., SAJANTILA, A., TAGLIABRACCI, A., GETHER, U., WERGE, T., RIVADENEIRA, F., HOFMAN, A., UITTERLINDEN, A. G., GIEGER, C., WICHMANN, H.-E., RÜTHER, A., SCHREIBER, S., BECKER, C., NÜRNBERG, P., NELSON, M. R., KRAWCZAK, M. and KAYSER, M. (2008). Correlation between genetic and geographic structure in Europe. *Curr. Biol.* **18** 1241–1248.
- LEE, S., WRIGHT, F. A. and ZOU, F. (2010). Control of population stratification by correlation-selected principal components. Preprint.
- LEUTENEGGER, A.-L., PRUM, B., GÉNIN, E., VERNY, C., LEMAINQUE, A., CLERGET-DARPOUX, F. and THOMPSON, E. A. (2003). Estimation of the inbreeding coefficient through use of genomic data. *Am. J. Hum. Genet.* **73** 516–523.
- LI, C. C. and HORVITZ, D. G. (1953). Some methods of estimating the inbreeding coefficient. *Am. J. Hum. Genet.* **5** 107–117.
- LIU, H., PRUGNOLLE, F., MANICA, A. and BALLOUX, F. (2006). A geographically explicit genetic model of worldwide human-settlement history. *Am. J. Hum. Genet.* **79** 230–237.
- MALÉCOT, G. (1969). *The Mathematics of Heredity*. Freeman, San Francisco, CA. [MR0258483](#)
- MARCHINI, J., CARDON, L. R., PHILLIPS, M. S. and DONNELLY, P. (2004a). The effects of human population structure on large genetic association studies. *Nat. Genet.* **36** 512–517.

- MARCHINI, J., CARDON, L. R., PHILLIPS, M. S. and DONNELLY, P. (2004b). Reply to “Genomic control to the extreme.” *Nat. Genet.* **36** 1129–1130; author reply 1131.
- MCCARTHY, M. L., ABECASIS, G. R., CARDON, L. R., GOLDSTEIN, D. B., LITTLE, J., IOANNIDIS, J. P. A. and HIRSCHHORN, J. N. (2008). Genome-wide association studies for complex traits: Consensus, uncertainty and challenges. *Nat. Rev. Genet.* **9** 356–369.
- MCKEIGUE, P. (2007). Population admixture and stratification in genetic epidemiology. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 1190–1213. Wiley, Chichester. [MR2391785](#)
- MCPPEEK, M. S. and SUN, L. (2000). Statistical tests for detection of misspecified relationships by use of genome-screen data. *Am. J. Hum. Genet.* **66** 1076–1094.
- MCVEAN, G. (2007). Linkage disequilibrium, recombination and selection. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 909–944. Wiley, Chichester. [MR2391785](#)
- MILLIGAN, B. G. (2003). Maximum-likelihood estimation of relatedness. *Genetics* **163** 1153–1167.
- MORRIS, A. and CARDON, L. (2007). Whole genome association. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 1238–1263. Wiley, Chichester. [MR2391785](#)
- NHGRI GWAS Catalog (2009). A catalog of published genome-wide association studies. Available at <http://www.genome.gov/gwastudies>.
- NOVEMBRE, J. and STEPHENS, M. (2008). Interpreting principal component analyses of spatial population genetic variation. *Nat. Genet.* **40** 646–649.
- PATTERSON, N., PRICE, A. L. and REICH, D. (2006). Population structure and eigenanalysis. *PLoS Genet.* **2** e190.
- PRENTICE, R. and PYKE, R. (1979). Logistic disease incidence models and case-control studies. *Biometrika* **66** 403–411. [MR0556730](#)
- PRICE, A. L., PATTERSON, N. J., PLENGE, R. M., WEINBLATT, M. E., SHADICK, N. A. and REICH, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38** 904–909.
- PRITCHARD, J. K. and DONNELLY, P. (2001). Case-control studies of association in structured or admixed populations. *Theor. Popul. Biol.* **60** 227–237.
- PRITCHARD, J. K. and PRZEWORSKI, M. (2001). Linkage disequilibrium in humans: Models and data. *Am. J. Hum. Genet.* **69** 1–14.
- PRITCHARD, J. K. and ROSENBERG, N. A. (1999). Use of unlinked genetic markers to detect population stratification in association studies. *Am. J. Hum. Genet.* **65** 220–228.
- PURCELL, S., NEALE, B., TODD-BROWN, K., THOMAS, L., FERREIRA, M. A. R., BENDER, D., MALLER, J., SKLAR, P., DE BAKKER, P. I. W., DALY, M. J. and SHAM, P. C. (2007). PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81** 559–575.
- RAKOVSKI, C. S. and STRAM, D. O. (2009). A kinship-based modification of the armitage trend test to address hidden population structure and small differential genotyping errors. *PLoS ONE* **4** e5825.
- RITLAND, K. (1996). Estimators for pairwise relatedness and individual inbreeding coefficients. *Genetical Research* **67** 175–185.
- ROBINSON, G. K. (1991). That BLUP is a good thing: The estimation of random effects. *Statist. Sci.* **6** 15–32. [MR1108815](#)
- ROSENBERG, N. A. and NORDBORG, M. (2006). A general population-genetic model for the production by population structure of spurious genotype-phenotype associations in discrete, admixed or spatially distributed populations. *Genetics* **173** 1665–1678.
- ROUSSET, F. (2002). Inbreeding and relatedness coefficients: What do they measure? *Heredity* **88** 371–380.
- SEAMAN, S. R. and RICHARDSON, S. (2004). Equivalence of prospective and retrospective models in the Bayesian analysis of case-control studies. *Biometrika* **91** 15–25. [MR2050457](#)
- SETAKIS, E., STIRNADEL, H. and BALDING, D. J. (2006). Logistic regression protects against population structure in genetic association studies. *Genome Res.* **16** 290–296.
- SLATKIN, M. (2002). The age of alleles. In *Modern Developments in Theoretical Population Genetics*, 3rd ed. (M. Slatkin and M. Veuille, eds.) 233–258. Oxford Univ. Press.
- SPIELMAN, R. S., MCGINNIS, R. E. and EWENS, W. J. (1993). Transmission test for linkage disequilibrium: The insulin gene region and insulin-dependent diabetes mellitus (iddm). *Am. J. Hum. Genet.* **52** 506–516.
- The Wellcome Trust Case Control Consortium (2007). Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447** 661–678.
- THOMPSON, E. A. (1975). The estimation of pairwise relationships. *Ann. Hum. Genet.* **39** 173–188. [MR0406572](#)
- THOMPSON, E. A. (1985). *Pedigree Analysis in Human Genetics*. Johns Hopkins Univ. Press, Baltimore, MD.
- THOMPSON, E. A. (2007). Linkage analysis. In *Handbook of Statistical Genetics*, 3rd ed. (D. J. Balding, M. Bishop and C. Cannings, eds.) **2** 1141–1167. Wiley, Chichester. [MR2391785](#)
- TIWARI, H. K., BARNHOLTZ-SLOAN, J., WINEINGER, N., PADILLA, M. A., VAUGHAN, L. K. and ALLISON, D. B. (2008). Review and evaluation of methods correcting for population stratification with a focus on underlying statistical principles. *Hum. Hered.* **66** 67–86.
- VOIGHT, B. F. and PRITCHARD, J. K. (2005). Confounding from cryptic relatedness in case-control association studies. *PLoS Genet.* **1** e32.
- WANG, Y., LOCALIO, R. and REBBECK, T. R. (2004). Evaluating bias due to population stratification in case-control association studies of admixed populations. *Genet. Epidemiol.* **27** 14–20.
- WANG, Y., LOCALIO, R. and REBBECK, T. R. (2005). Bias correction with a single null marker for population stratification in candidate gene association studies. *Hum. Hered.* **59** 165–175.
- WEALE, M. E., WEISS, D. A., JAGER, R. F., BRADMAN, N. and THOMAS, M. G. (2002). Y chromosome evidence for Anglo-Saxon mass migration. *Mol. Biol. Evol.* **19** 1008–1021.

- WEINBERG, C. R. (1999). Methods for detection of parent-of-origin effects in genetic studies of case-parents triads. *Am. J. Hum. Genet.* **65** 229–235.
- WEIR, B. S., ANDERSON, A. D. and HEPLER, A. B. (2006). Genetic relatedness analysis: Modern data and new challenges. *Nat. Rev. Genet.* **7** 771–780.
- YU, J., PRESSOIR, G., BRIGGS, W. H., BI, I. V., YAMASAKI, M., DOEBLEY, J. F., MCMULLEN, M. D., GAUT, B. S., NIELSEN, D. M., HOLLAND, J. B., KRESOVICH, S. and BUCKLER, E. S. (2006). A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat. Genet.* **38** 203–208.
- ZHANG, S., ZHU, X. and ZHAO, H. (2003). On a semiparametric test to detect associations between quantitative traits and candidate genes using unrelated individuals. *Genet. Epidemiol.* **24** 44–56.
- ZHAO, K., ARANZANA, M. J., KIM, S., LISTER, C., SHINDO, C., TANG, C., TOOMAJIAN, C., ZHENG, H., DEAN, C., MARJORAM, P. and NORDBORG, M. (2007). An arabidopsis example of association mapping in structured samples. *PLoS Genet.* **3** e4.
- ZHENG, G., FREIDLIN, B. and GASTWIRTH, J. L. (2006). Robust genomic control for association studies. *Am. J. Hum. Genet.* **78** 350–356.
- ZHENG, G., FREIDLIN, B., LI, Z. and GASTWIRTH, J. L. (2005). Genomic control for association studies under various genetic models. *Biometrics* **61** 186–192. [MR2135859](#)
- ZHENG, G., LI, Z., GAIL, M. H. and GASTWIRTH, J. L. (2010). Impact of population substructure on trend tests for genetic case-control association studies. *Biometrics*. To appear.