

Quality Control Procedures for Genome-Wide Association Studies

UNIT 1.19

Stephen Turner,¹ Loren L. Armstrong,² Yuki Bradford,¹ Christopher S. Carlson,³ Dana C. Crawford,¹ Andrew T. Crenshaw,⁴ Mariza de Andrade,⁵ Kimberly F. Doheny,⁶ Jonathan L. Haines,¹ Geoffrey Hayes,² Gail Jarvik,⁷ Lan Jiang,¹ Iftikhar J. Kullo,⁸ Rongling Li,⁹ Hua Ling,⁶ Teri A. Manolio,⁹ Martha Matsumoto,⁵ Catherine A. McCarty,¹⁰ Andrew N. McDavid,³ Daniel B. Mirel,⁴ Justin E. Paschall,¹¹ Elizabeth W. Pugh,⁶ Luke V. Rasmussen,¹⁰ Russell A. Wilke,¹² Rebecca L. Zuvich,¹ and Marylyn D. Ritchie¹

¹Center for Human Genetics Research, Department of Molecular Physiology & Biophysics, Vanderbilt University, Nashville, Tennessee

²Division of Endocrinology, Metabolism, and Molecular Medicine, Feinberg School of Medicine, Northwestern University, Chicago, Illinois

³Cancer Prevention, Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, Washington

⁴Genetic Analysis Platform and Program in Medical and Population Genetics, Broad Institute, Cambridge, Massachusetts

⁵Division of Biostatistics and Informatics, Department of Health Sciences Research, Mayo Clinic College of Medicine, Rochester, Minnesota

⁶Center for Inherited Disease Research, Johns Hopkins University, Baltimore, Maryland

⁷Department of Genome Sciences, University of Washington, Seattle, Washington

⁸Division of Cardiovascular Diseases, Department of Medicine, Mayo Clinic, Rochester, Minnesota

⁹Office of Population Genomics, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland

¹⁰Biomedical Informatics Research Center, Marshfield Clinic Research Foundation, Marshfield, Wisconsin

¹¹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland

¹²Division of Clinical Pharmacology, Department of Medicine, Vanderbilt University, Nashville, Tennessee

ABSTRACT

Genome-wide association studies (GWAS) are being conducted at an unprecedented rate in population-based cohorts and have increased our understanding of the pathophysiology of complex disease. Regardless of context, the practical utility of this information will ultimately depend upon the quality of the original data. Quality control (QC) procedures for GWAS are computationally intensive, operationally challenging, and constantly evolving. Here we enumerate some of the challenges in QC of GWAS data and describe the approaches that the electronic MEDical Records and Genomics (eMERGE) network is using for quality assurance in GWAS data, thereby minimizing potential bias and error in GWAS results. We discuss common issues associated with QC of GWAS data, including data file formats, software packages for data manipulation and analysis, sex chromosome anomalies, sample identity, sample relatedness, population substructure, batch effects, and marker quality. We propose best practices and discuss areas of ongoing and future research. *Curr. Protoc. Hum. Genet.* 68:1.19.1-1.19.18. © 2011 by John Wiley & Sons, Inc.

Keywords: genome-wide association studies • GWAS • quality control • QC • biobanks • electronic medical records • eMERGE

Genetic
Mapping

1.19.1

Supplement 68

INTRODUCTION

Genome-wide association studies (GWAS) are commonly used to identify common single nucleotide polymorphisms (SNPs) that influence human traits. GWAS have been conducted at increasing frequency using a case-control, population-based perspective, with cross-sectional study designs (Klein et al., 2005; Frayling, 2007; Willer et al., 2008; Hindorff et al., 2009; Kathiresan, 2009; Newton-Cheh et al., 2009). More recently, GWAS are being conducted in cohorts that are clinic-based (Link et al., 2008; Daly et al., 2009; Thompson et al., 2009; Barber et al., 2010). As a result, GWAS may soon move the field of genomics into clinical practice.

Whether the goal is to identify predictors of outcomes or to discover new biology underlying a trait of interest, the capability of GWAS to identify true genetic associations depends upon the overall quality of the data. Even simple statistical tests of association are compromised in the context of genome-wide SNP data that have not been properly cleaned up, potentially leading to false-negative and false-positive associations. Additionally, problems with the overall data quality will likely affect downstream analyses and studies beyond the initial GWAS. For example, the National Human Genome Research Institute (NHGRI) actively maintains an online catalog of GWAS results and associated publications (Hindorff et al., 2009), which stimulates downstream studies of replication and characterization in independent populations. Compromised data quality in the discovery phase may lead to false-positive results that are carried forward into replication studies at great cost in terms of both time and expense. Also, the National Institutes of Health (NIH) now mandates that secure, encrypted copies of primary GWAS data funded by NIH be made publicly available (with controlled access) for secondary analyses. These accessible datasets are maintained by the National Center for Biotechnology Information (NCBI) in the database of Genotypes and Phenotypes (dbGaP). dbGaP provides both open and controlled access, which allow for both broad release of nonsensitive information and restricted access to datasets involving genomic data and phenotypic information, respectively (Mailman et al., 2007). Data access through dbGaP is commonly used for replication and meta-analysis, both of which will be compromised by poor-quality data.

Genotyping technology and allele-calling algorithms continue to improve, and quality-

improvement strategies continue to ensure that only reliable, rigorously scrutinized markers and samples are used for analysis. Reconciling genetic data with clinical and self-reported data (e.g., sex or familial relationships) can potentially identify sample-identity problems caused by sample-handling mishaps. Batch effects, population stratification, and sample relatedness can confound genetic association analyses and can lead to excessive type I and type II errors. Here, we discuss methods that can be used to detect and account for various data-quality issues to better ensure the integrity of the primary GWAS as well as its downstream applications.

The eMERGE (electronic Medical Records and GENomics) Network (<https://www.gwas.net>) is an NHGRI-supported consortium of five institutions charged with exploring the utility of DNA repositories coupled to Electronic Medical Record (EMR) systems for advancing discovery in genome science (McCarty et al., 2010). Genome-wide genotyping has been performed on ~17,000 samples across the eMERGE network at the Broad Institute and at the Center for Inherited Disease Research (CIDR) using the Illumina 660W-Quad or 1 M-Duo Beadchips. Each study site conducts a GWAS, in addition to a number of cross-network analyses. These studies adhere to NIH's data-sharing policies, and all data generated in this study will be available on dbGaP (Mailman et al., 2007). Due to the complexity involved in a single-site GWAS, in addition to the combining of data and results across study sites, it became clear that a unified QC pipeline was imperative.

Others have discussed quality control procedures for genotypic data (Broman, 1999; Chanock et al., 2007; Miyagawa et al., 2008; Laurie et al., 2010). The goal of this manuscript is to provide a tutorial to instruct investigators on QC procedures that should be performed prior to GWAS data analysis. The procedures discussed here were developed by the genomics group of the eMERGE network, where phenotyping and other sample information is obtained through sophisticated mining of the EMR. This protocol can be applied to many GWAS studies, regardless of phenotyping strategy. Given that most of the genotyping data available for GWAS are currently SNP-based, we will limit our discussion to these biallelic markers, and QC procedures for CNV analysis will not be discussed here. Figure 1.19.1 shows a flowchart overview of the entire QC process, where each step is discussed in detail in the following sections.

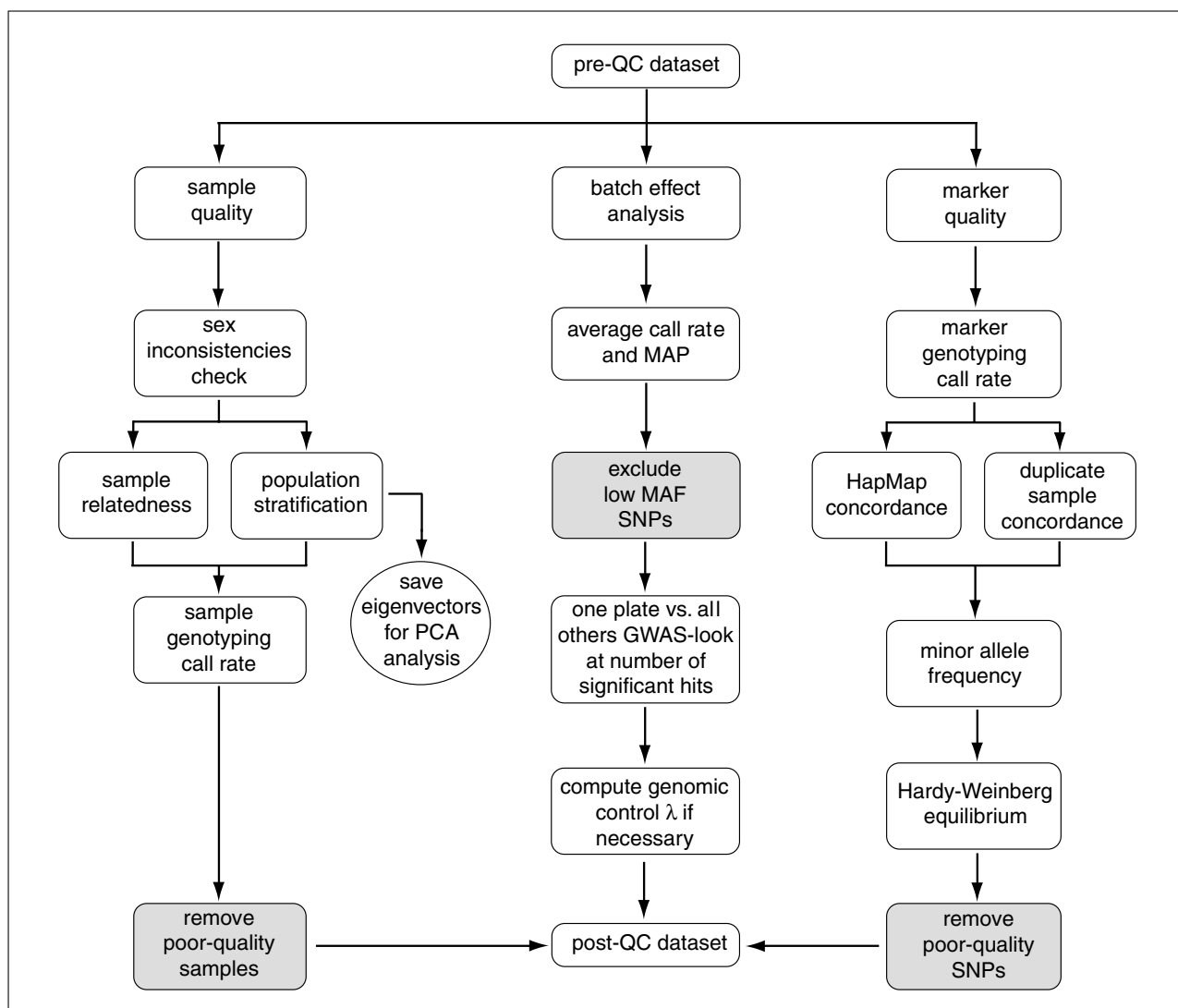


Figure 1.19.1 A flowchart overview of the entire GWAS QC process. Each topic is discussed in detail in the corresponding section in the text. Squares represent steps, ovals represent input or output data, and trapezoids represent filtering of data.

GWAS DATA FORMAT

Regardless of the underlying study design (such as family-based or population-based), the most commonly used format for genetic data is the linkage, or pedigree file format (pedfile). This file contains one individual per row, where the first six columns are identifying information (family ID, individual ID, father ID, mother ID, sex, phenotype), and the remaining columns are genotypes (two columns per genotype; one for each allele). The genotype column pairs correspond to an ordered set of SNP markers present in an associated file (.map or .bim). Additional phenotypes can also be stored in separate files consisting of family ID, individual ID, then extra columns representing additional phenotypes. There are several variations on pedfile format, including transposed (long) formats (tped), and com-

pressed (binary) formats. Descriptions of these file formats can be found on the PLINK homepage (Table 1.19.1). PLINK is a freely available, open-source, cross-platform application for QC and analysis of GWAS data (Purcell et al., 2007). We used PLINK for implementing most of the eMERGE network's QC pipeline.

An important issue when creating a pedfile for QC analysis is the choice of strand orientation to use for allele calls (i.e., forward or reverse complement). While forward strand is a commonly used allele-coding scheme, Illumina has developed a consistent and simple method to ensure uniformity in genotype call reporting that uses the polymorphism itself and the contextual surrounding sequence ("TOP/BOT" strand and "A/B" allele coding; Illumina Technical Note; see Internet

Table 1.19.1 Useful Software Packages for Data Management, Quality Control, and Statistical Analysis in Genome-Wide Association Studies

Software package	URL	Purpose
PLINK	http://pngu.mgh.harvard.edu/~purcell/plink/	Free, open-source GWAS analysis software package. Contains many tools for data management, quality control, and statistical analysis. (PC, Mac, Linux).
PLATO	https://chgr.mc.vanderbilt.edu/plato	<i>PL</i> atform for the Analysis, Translation, and Organization of large-scale data. Software for GWAS analysis similar to PLINK.
R	http://www.r-project.org/	Free, open-source statistical computing software with excellent graphical capabilities (PC, Mac, Linux)
Eigensoft	http://genepath.med.harvard.edu/~reich/Software.htm	Free, open-source software for performing principal components analysis based method for detecting and correcting for population stratification in GWAS (Linux only)
Structure	http://pritch.bsd.uchicago.edu/structure.html	Free, open-source software for inferring the presence of distinct populations and assigning individuals to those populations for a stratified analysis (Windows, Mac, Linux)
MySQL Workbench	http://wb.mysql.com/	Free, open-source software for creating, administering, and querying relational databases. This is helpful for subsetting data, merging results, and joining QC metrics (e.g., HWE) to final association results (Windows, Mac, Linux).

Resources). Since 2005, the database of genetic variation (dbSNP; Sherry et al., 2001) has used this designation for all SNP entries. We used “TOP/BOT” strand orientation for eMERGE. Choice of strand orientation might depend on the strand orientation of other data used in a combined analysis or of a reference set used for imputation. The goal is to ensure uniformity in genotype call reporting, which is critically important in downstream analyses, reporting, and annotation.

SAMPLE QUALITY

Sex inconsistencies and chromosomal anomalies

One of the first procedures that should be implemented in any GWAS QC protocol is checking for potential sample-identity problems that typically result from sample-handling errors. One of the easiest ways to discover potential sample-handling issues that result in mix-ups is by checking the reported sex of each individual against that predicted by the genetic data. The `—check-sex` option in PLINK uses X chromosome heterozygosity rates to determine sex empirically, then reports individuals for whom the sex recorded in the pedfile does not match the predicted sex based on genetic data (example output and ex-

planation shown in Table 1.19.2). If discrepancies are found (e.g., an individual is recorded being female but appears homozygous for every X chromosome marker), the EMR or any available study questionnaires should be reviewed to make a determination whether there was a sample-handling mistake that caused a sample mix-up. Checking X chromosome heterozygosity may also reveal sex chromosome anomalies such as Turner syndrome (females having karyotype XO), Klinefelter syndrome (males having karyotype XXY), mosaic individuals (XX/XO, XX/XXY), or females with large stretches of loss-of-heterozygosity on the X chromosome who are otherwise phenotypically normal. In one eMERGE cohort, the rate of XX/XO mosaicism was 0.08%. The rate of XXY was 0.025%, and the rate of XXY/XY mosaicism was 0.05%.

X chromosome heterozygosity is a fairly sensitive heuristic to detect sample swaps, but not very specific. A variety of factors besides a crude sample mix-up will affect heterozygosity. Furthermore, if the goal is to enumerate as many samples with atypical sex karyotypes as possible, then X heterozygosity alone will not detect abnormalities such as triple X or XYY or homozygous X Klinefelter syndrome. Examining the intensity of probe binding on the sex chromosomes will better resolve these

Table 1.19.2 Example Table Showing Output from `—check-sex` Routine Using PLINK^a

IID	PEDSEX	SNPSEX	STATUS	F	Explanation
1	1	1	OK	0.98	Male
2	2	2	OK	0.03	Female
3	2	1	PROBLEM	0.99	Recorded female, genetically male
4	1	2	PROBLEM	0.02	Recorded male, genetically female
5	2	0	PROBLEM	0.28	Likely a female with sex chromosome anomaly (e.g., XX/XO mosaic, loss-of-heterozygosity on X)
6	1	0	PROBLEM	0.35	Likely a male with sex chromosome anomaly (e.g., XXY or XX/XY mosaic)

^aIID = individual id; PEDSEX = sex as recorded in pedfile (1 = male, 2 = female); SNPSEX = sex as predicted based on genetic data (1 = male, 2 = female, 0 = unknown).

cases. Illumina calls this intensity LogR ratio. On Affymetrix systems, it is simply known as probe intensity. These metrics, once suitably normalized, are roughly linear in copy number. Because there are tens of thousands of loci on the X chromosome on modern platforms, it is appropriate to examine a subsample of markers and then take a measure of central tendency of each sample such as the median or mean intensity. The intensity plot provides a visualization of the intensity of X and Y probes (Fig. 1.19.2). It is expected that females should have low Y intensity and high X intensity (bottom right corner), and the males should show similar levels of X and Y intensities (top left corner). We also observed two individuals with mislabeled sex as well as several individuals with XXY. Structural chromosomal variation can be identified using intensity-only probes to calculate loss of heterozygosity and abnormal copy numbers using B allele frequency and LogR Ratio plots (Fig. 1.19.3). The B allele frequency plot is the amount of B allele observed in a probe that should concentrate at zero for zero copy, at 0.5 for one copy, and at 1 for two copies. LogR ratio indirectly measures copy number of each SNP by plotting the ratio of observed to expected hybridization intensity (Simon-Sanchez et al., 2007). It is expected to concentrate at zero; sometimes an upward or downward bump is observed meaning amplification or deletion, respectively. These two plots can be obtained using Illumina BeadStudio/GenomeStudio.

Depending on the aims of the study, often these individuals are not eliminated from the study due to sex chromosome anomalies alone. Even in carefully collected samples, the numbers of samples with discrepant self-reported sex having a normal karyotype is appreciable, and sample-processing pipelines need to have these checks in place to detect such potential sample swaps. While it may be possible to go back to the original data to reconcile chromosomal anomalies found using genetic data, researchers need to be aware of any ethical issues that may arise concerning return of results, and these issues should be considered and resolved prior to revisiting the EMR.

Sample relatedness

Another way to simultaneously examine both sample identity and pedigree integrity is by reconciling genomic data with self-reported relationships between individuals (if available). Although this has consequences that impact which analytical approaches are appropriate for the downstream association study, having related samples in the dataset makes it possible to further investigate potential DNA sample mix-ups. Using dense marker data obtained in GWAS, it is easy to compute pairwise kinship estimates between every individual in the study using the `—genome` option in PLINK. This procedure need not be performed on the entire GWAS dataset; using only 100,000 markers will also yield stable estimates of kinship coefficients.

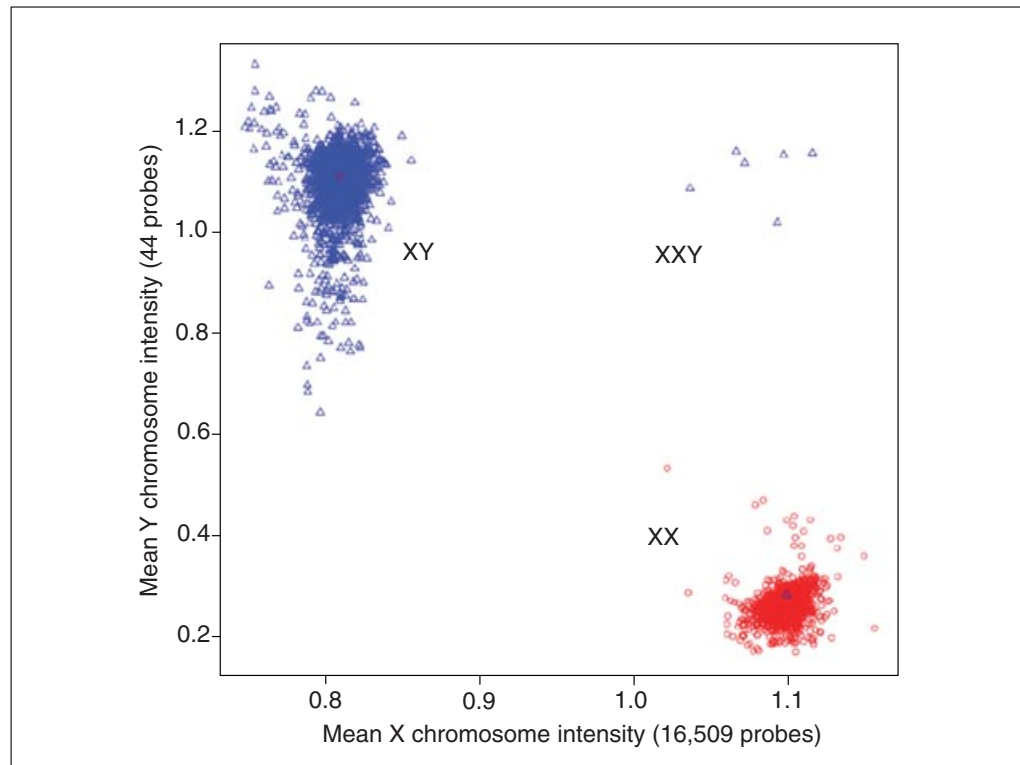


Figure 1.19.2 Visualization of X and Y probe intensities. The x-axis and y-axis represent the sum of the average over all probes for the normalized Cartesian intensity for allele A and the average over all probes for the normalized Cartesian intensity for allele B using all probes available on the X chromosome and Y chromosome, respectively. The XX (female, red circles) and XY (male, blue triangles) subjects are shown on the bottom right corner and on the top left corner, respectively. The plot reveals two mislabeled individuals (one male with the female cluster, and one female with the male cluster). Several XXY individuals are also clearly visible (upper right corner). For color version of this figure go to <http://www.currentprotocols.com/protocol/hg0119>.

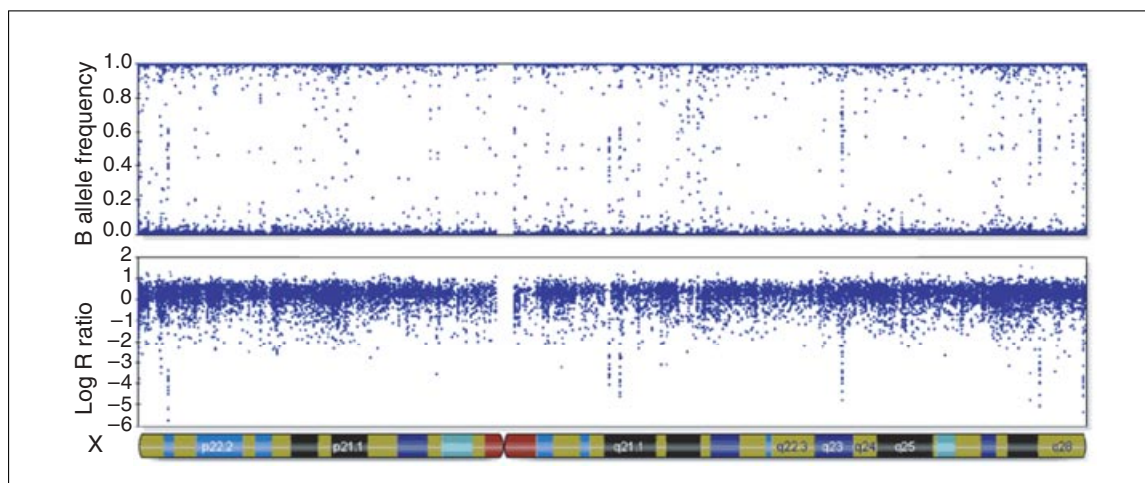


Figure 1.19.3 Copy Number and allelic variation to detect anomalies on the X chromosome. The top plot shows the B-Allele frequencies for all probes for one sample with total loss of heterozygosity (LOH) on the X chromosome. The bottom plot shows the copy number variation from the same sample on X chromosome. Both plots are helpful to detect regions of LOH and/or copy number variation such as deletion and amplification.

In addition to reporting the relationship type as reported using pedigree data (e.g., siblings, parent-child, unrelated), this procedure will also calculate the proportion of loci where two individuals share zero, one, or two alleles identical by descent (IBD). Individuals sharing two alleles IBD at every locus are monozygotic twins, or else the pair is actually a single sample processed twice. Individuals sharing zero alleles IBD at every locus are unrelated. Individuals sharing one allele IBD at every locus are parent-child pairs. On average, siblings share zero, one, and two alleles IBD at 25%, 50%, and 25% of the genome, respectively. Using these data, the proportion of loci sharing one allele IBD (Z1) can be plotted by the proportion of loci where individuals share zero alleles IBD (Z0) and points color coded by the relationship type. For clarity, this plot can be restricted to points where the overall kinship coefficient is ≥ 0.05 , as most of the individuals where kinship ≤ 0.05 will be unrelated. This will produce a plot as shown in Figure 1.19.4. Detailed instructions on producing this graphic using R (R Development Core Team; see Internet Resources) can be found online (S.D. Turner; see Internet Resources).

If it is believed that pedigree records obtained through the original data are accurate, then a point out of place (e.g., points colored as unrelated showing up where most of the parent-offspring pairs cluster) would be

indicative of either nonpaternity, adoption, sample mix-up, or duplicate processing of a single individual. Further investigation employing the original data can be used to attempt to identify the problem. It is also worth noting in studies where datasets from multiple sites are combined that it is possible that the same participant is present in more than one study. These two data points would appear genetically identical across sites.

In addition to potentially discovering sample-handling issues, visualizing sample relatedness as shown in Figure 1.19.4 also reveals any cryptic relatedness that may be present in the study sample. Figure 1.19.4 shows that many individuals who indicated that they were unrelated (black points) or distantly related (blue points) line up along the diagonal in this plot. These individuals represent second-, third-, fourth-, and fifth-degree relatives. If treated as independent samples in the downstream analyses, having many related samples in the dataset would result in increased type I and type II errors; thus, analytical methods such as mixed-model regression (Aulchenko et al., 2007) must be used in place of simple linear or logistic regression. Figure 1.19.5 shows another way to visualize the degree of relatedness by plotting a histogram of the distribution of kinship coefficients over 0.05 between all pairs of individuals in the dataset.

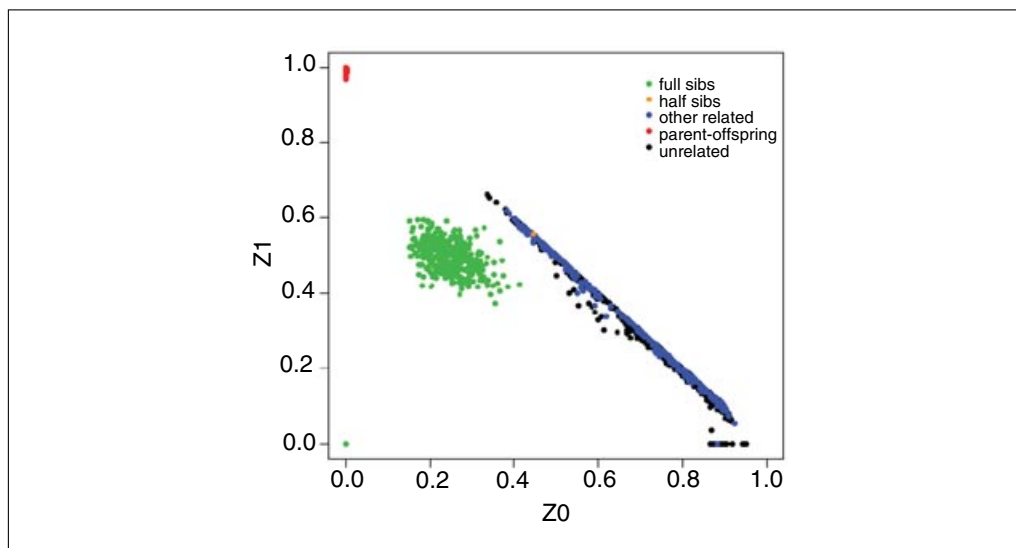


Figure 1.19.4 Points in this plot show pairs of individuals plotted by their degree of relatedness: the proportion of loci where the pair shares one allele IBD (Z1) by the proportion of loci where the pair shares zero alleles IBD (Z0). These values are obtained from PLINK using the `—genome` option. Pairs are color-coded by the type of relationship determined by the pedigree information embedded in the pedfile (also reported by PLINK). This plot omits pairs of individuals having an overall kinship coefficient ≥ 0.05 for clarity. There is a pair of monozygotic twins represented by a point in the lower left at (0,0), because they share two alleles IBD at every locus across the genome. For color version of this figure go to <http://www.currentprotocols.com/protocol/hg0119>.

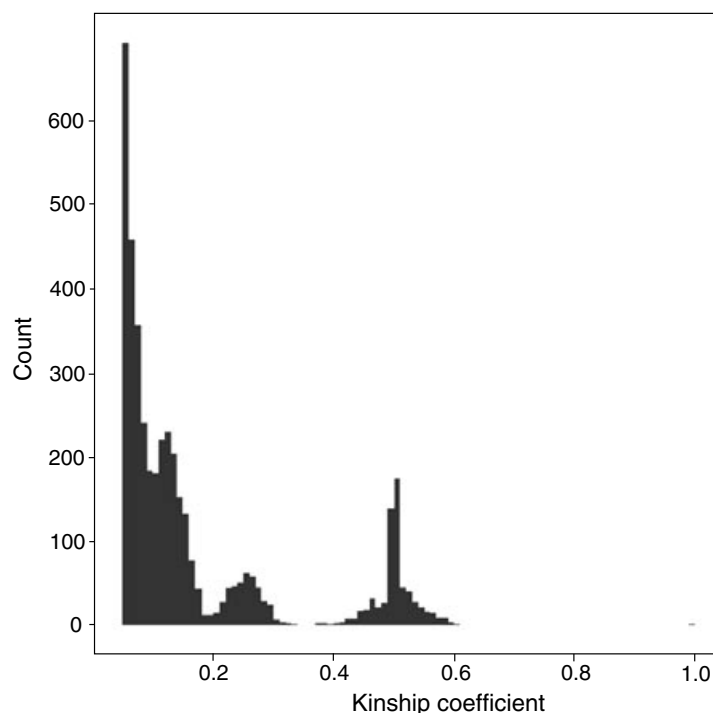


Figure 1.19.5 Histogram showing the distribution of pairwise kinship coefficients (where kinship coefficient is greater than 0.05). The peak over 0.5 represents first degree relatives (parent-offspring, full siblings). The peak over 0.25 represents second-degree relatives (half siblings, avuncular, grandparent-grandchild). Third- and fourth-degree relatives begin to blend into more distantly related samples between zero and 0.125.

Population substructure

Population stratification occurs when the study samples comprise multiple groups of individuals who differ systematically in both genetic ancestry and the phenotype under investigation. Spurious apparent associations would be due to differences in ancestry rather than true association of alleles to disease (Cardon and Palmer, 2003). Thus, it is critical to check for population stratification within the study samples and leverage this information to inform the downstream analyses.

One strategy for avoiding bias induced by population stratification is to ensure that study samples are drawn from a relatively homogenous population. One of the sites in the eMERGE network represents such a sample, as over 98% of the study sample self-reported “Caucasian” on a study questionnaire. This percentage is consistent with data from the 2000 Census (see Internet Resources), and self-reporting often shows very high correspondence with genetically inferred ancestry (Tang et al., 2005). Some clinics record ethnicity via observer report (typically a clerk or nurse’s aide). Even in this setting, observer-reported ancestry closely matches genetically inferred

ancestry, especially for populations of European descent (Dumitrescu et al., 2010). However, population-based diverse samples are often desirable for genetic association studies focused on characterizing previous GWAS or candidate gene discoveries made in one population (Manolio, 2009). Furthermore, combining samples from multiple sites for a joint analysis may result in population stratification in the combined sample, if both allele frequencies and outcomes differ between sites.

Statistical methodology has been developed and implemented into software to aid in detecting and adjusting for population stratification in GWAS. Genomic control (Devlin and Roeder, 1999; Reich and Goldstein, 2001) aims to control for population stratification by first estimating an inflation factor (λ), then adjusting all of the test statistics downward by this factor (Fig. 1.19.1). Several variations on genomic control have been developed, and a recent review and critical evaluation of genomic control methods (Dadd et al., 2009) recommended genomic control F (GCF; Devlin et al., 2004) as the most appropriate variation. GCF does not assume that the inflation factor is measured without error, and

refines this factor accordingly. Structured association (Pritchard et al., 2000), implemented in the STRUCTURE software (see Internet Resources), uses genotype data to infer population structure and subsequently performs tests of association within each inferred subpopulation. STRUCTURE may also be used to identify individual samples that do not cluster with the majority of the samples. These samples can then be eliminated from the analysis.

Because the risk of confounding by population stratification may increase with sample size (i.e., confounded results become more significant with larger samples; Marchini et al., 2004), and because large sample GWAS are becoming increasingly common, another method has been developed that utilizes large samples and thousands of markers throughout the genome to adjust for population structure. Eigenstrat analysis (Patterson et al., 2006; Price et al., 2006) uses principal components analysis to explicitly detect and adjust for population stratification on a genome-wide scale with large sample sizes in a computationally

efficient manner. This method may be preferred over a stratified analysis because the combined sample often yields more powerful statistical tests, even after adjusting for significant eigenvectors (Zhang et al., 2008).

Eigensoft is freely available open-source software for conducting Eigenstrat analyses, available online (Table 1.19.1). Running Eigensoft requires dense genotyping coverage. We recommend using all the default options, including 100,000 randomly chosen high-quality markers. There are several SNPs in the HLA region on chromosome 6, in the lactase locus on chromosome 2, and in the inversion regions on 8p23 and 17q21.31 common in populations of European ancestry (Novembre et al., 2008), which are sources of stratification that will often appear in the top principal components. While one may exclude these SNPs from such an analysis, it is unknown if any similar inversions exist at appreciable frequency in non-European populations. Thus, it may be preferable to detect and correctly interpret the analysis

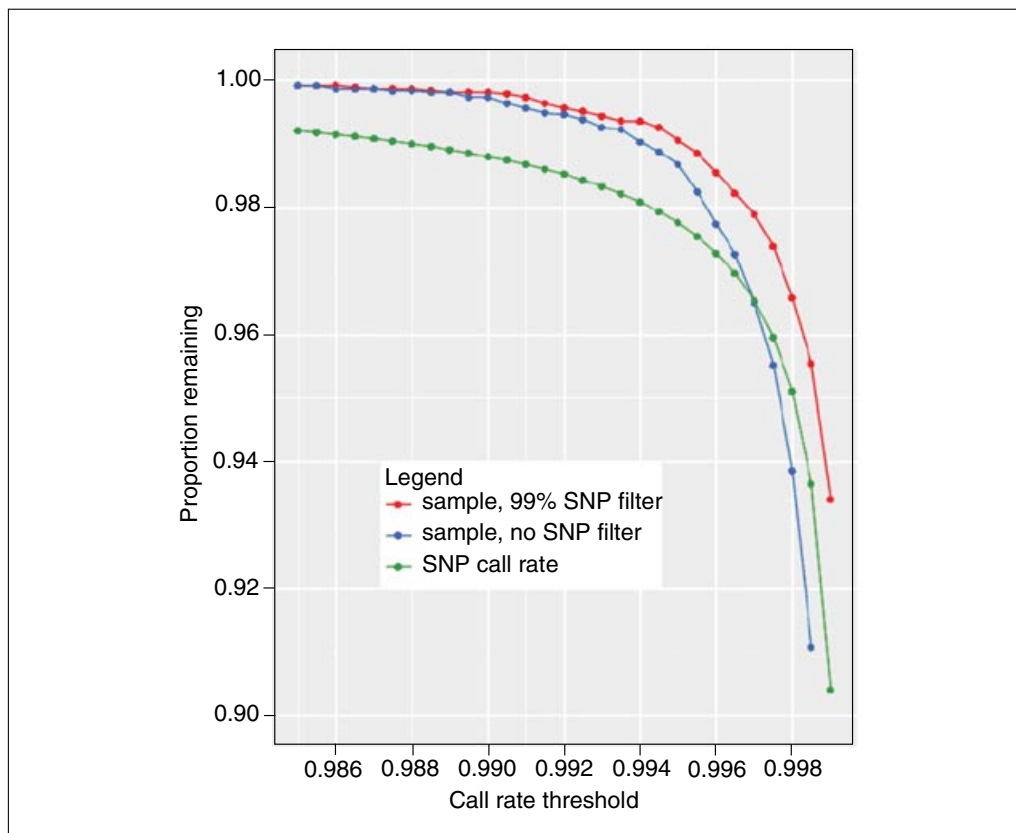


Figure 1.19.6 Proportion of SNPs or samples remaining as call rate threshold increases. The green line shows the proportion of SNPs remaining when SNPs are discarded if they fall below the given genotyping efficiency threshold. The blue line shows the proportion of samples remaining, while the red line shows the proportion of samples remaining if a 99% call rate threshold is applied to eliminate poor-quality markers first. For color version of this figure go to <http://www.currentprotocols.com/protocol/hg0119>.

including these regions rather than avoiding specific regions.

The Eigensoft analysis will result in the computation of 10 principal components. If any of these eigenvectors are significantly associated with the phenotype under study, it is recommended that these eigenvectors be adjusted for in any downstream analysis to correct for any bias due to population stratification. Alternatively, if it is expected that only a very small number of samples represent ethnic outliers in the study population, using Eigenstrat with iterative outlier removal and reconciling these individuals with other ancestry information such as self-report could be used to identify a coherent set of ethnic outliers (which are identified as outliers *and* self-reported as outliers) to potentially exclude from the analysis, rather than adjusting for many eigenvectors during the analysis only to retain a very small number of samples.

Sample genotyping efficiency/call rate

Genotyping efficiency, or call rate, is an issue which will be discussed in greater depth in the “Marker Quality” section below. A large proportion of SNP assays failing on an individual DNA sample may be indicative of a poor-quality DNA sample, which could lead to aberrant genotype calling. Samples with low genotyping efficiency, or call rate, should be eliminated from further analysis. A recommended threshold is 98% to 99% efficiency, after first removing markers which have a low genotype call rate across samples. The suggested 98% to 99% threshold is an approximate threshold—the exact threshold may vary from study to study depending on the genotyping platform used, quality of the DNA samples used, and the variability in human and equipment error in genotyping. The threshold should be determined based on a goal whereby a balance minimizing the number of samples dropped and maximizing genotyping efficiency is attained. Figure 1.19.6 shows the proportion of samples (red and blue lines) or SNPs (green line) remaining at different call rate thresholds. Genotyping efficiency can be checked using the `—missing` option in PLINK. This will produce a file showing genotype missingness rate (1-efficiency) for each individual (proportion of SNPs which failed on each sample), and for each SNP (proportion of individuals for which no genotype was called). Samples below a desired threshold can be eliminated from any downstream analyses by using the `—mind` option in PLINK. Geno-

typing efficiency is also an important marker QC step, and is discussed below.

MARKER QUALITY

Marker genotyping efficiency/call rate

As mentioned in the “Sample genotyping efficiency” section above, marker genotyping efficiency (the proportion of samples with a genotype call for each marker) is a good indicator of marker quality. SNP assays that failed on a large number of samples are poor assays, and are likely to result in spurious data. A recommended threshold for removing SNPs with low call rate is approximately 98% to 99%, although as mentioned in the “Sample genotyping efficiency” section, this threshold may vary from study to study. Marker genotyping efficiency can be reviewed using the `—missing` option in PLINK. We recommend removing poor-quality SNPs before running the sample genotyping efficiency check discussed above, so that fewer samples will be dropped from the analysis simply because they were genotyped with SNP assays that had poor performance (see Fig. 1.19.6). Markers can be removed based on call rate by using the `—geno` option, followed by a threshold for a lower limit of missingness (e.g., `—geno 0.02` would remove SNPs with more than 2% missing, i.e., less than a 98% call rate).

Control sample reproducibility/HapMap concordance

It is advantageous to incorporate internal controls in the genotyping pipeline to estimate genotyping reproducibility rate and for selecting which markers to eliminate based on poor reproducibility. Many studies routinely genotype DNA samples from the HapMap cell lines (International HapMap Consortium, 2003, 2007). In addition to providing samples of known ancestry to anchor the STRUCTURE analysis discussed in the “Population substructure” section above, genotype calls on HapMap samples can be compared to the corresponding publicly available reference genotypes to estimate the degree of concordance. Genotyping for two centers in the eMERGE network was performed by CIDR, which considered any SNP having more than one replicate error on HapMap samples run with the study samples to be a technical failure, and only intensity data were released for these markers. CIDR also considered SNPs technical failures if the SNP had a call rate <85%, if the absolute difference in call rate between

sexes is greater than 2.5%, if the absolute difference in heterozygosity between sexes is greater than 7%, or if cluster separation <0.20 . Samples at three other centers in the eMERGE network were genotyped at the Broad Institute, where technical failure was determined by call rate 95%, GenTrain score <0.6 (a statistical measure from Illumina's clustering algorithm; Illumina GeneCall Data Analysis Software; see Internet Resources), cluster separation <0.4 , or more than one replicate error. It is also advantageous to build in duplicate samples to estimate the reproducibility rate within genotyping batches. By design, both CIDR and Broad include HapMap control samples and duplicate samples across all plates in the study.

It is anticipated that, for accurate genotyping data, duplicate reproducibility and HapMap concordance of $>99\%$ would be expected. We removed any SNPs which had one or more discordant calls on duplicate samples. Both HapMap concordance and replicate sample concordance can be checked using the concordance procedure in the PLATO software (Grady et al., 2010; Table 1.19.1)

or by using the `—genome —rel-check` options in PLINK. Using HapMap trio samples, it is also possible to inspect each SNP for Mendelian inconsistencies, which indicate genotyping errors if pedigree information is correct. Mendelian inconsistencies can be assessed using the `—mendel` option in PLINK. PLINK only detects Mendelian inconsistencies in full trios. The Mendelian-error procedure in PLATO will also evaluate Mendelian consistency in sib-pairs or in parent-offspring pairs where a full trio is unavailable.

We recommend removing or flagging any SNPs that have one or more Mendelian errors on HapMap control samples. While it may be possible to look for Mendelian inconsistencies using study samples, removing these SNPs could potentially be filtering out a phenotype-specific copy-number variant. If this is the case, there will likely be more than three genotype clusters. The extra clusters or parts of them will be missing or miscalled. For instance, for a locus with alleles A and B, A, AA, and AAB may all cluster together unless the SNP is re-called with a specific model in mind.

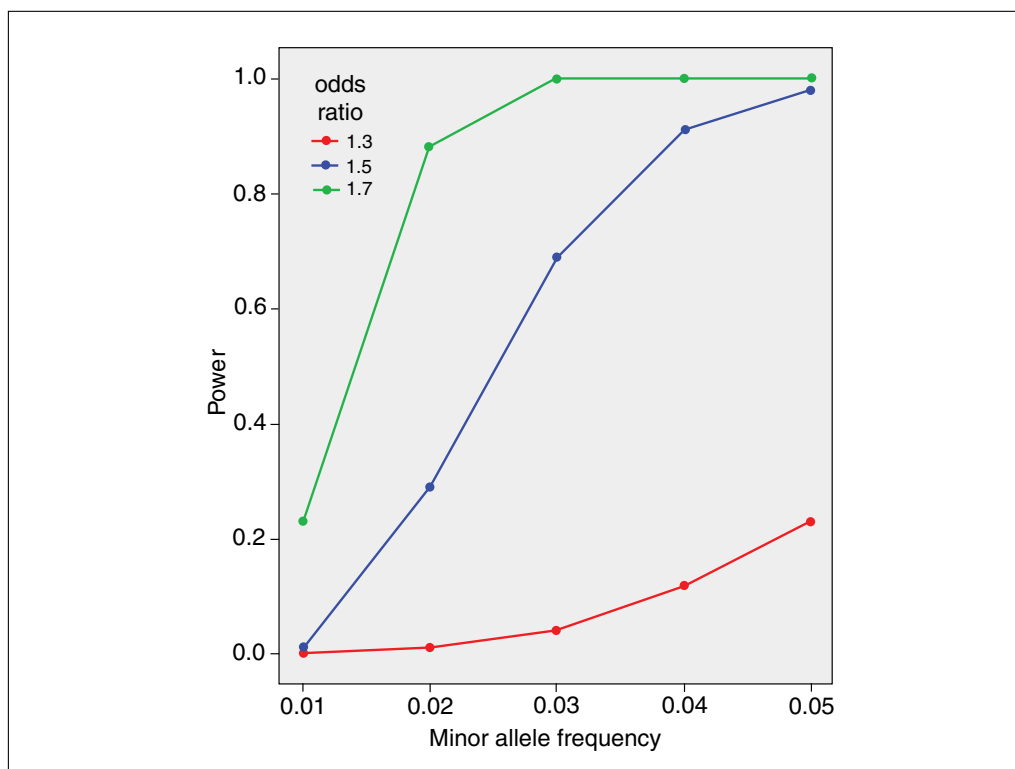


Figure 1.19.7 This shows the power to detect an association at genome-wide significance ($p < 5 \times 10^{-8}$), assuming the actual causal SNP is genotyped in a case-control study consisting of 5000 cases and 5000 controls of a common disease with 10% prevalence under an additive model at several different odds ratios. Note that when the MAF is low, power is extremely low even for very large effects (odds ratio = 1.7).

Minor allele frequency

It is also important to filter SNPs based on minor allele frequency, because statistical power is extremely low for rare SNPs. Figure 1.19.7 shows that the power to detect an association in a large dataset ($n = 10,000$) with a relatively large effect (odds ratio between 1.3 and 1.7) is extremely low for rare SNPs ($<1\%$ frequency). In addition to limited power for SNPs with low MAF, these SNPs also have the potential to lead to spurious associations due to either genotyping errors or population stratification. For SNPs with low MAF, the clustering algorithms for making genotype calls can be challenged. Therefore, it is conceivable that although many calls are made, and the marker call rate is acceptable, review of the cluster plots will show that the low MAF lead to poor clusters. Finally, if specific alleles are present only in certain ancestral populations, with low MAF in those populations, they can lead to associations that are due to potential stratification issues rather than disease association. We recommend removing any ex-

tremely rare SNPs (including any monomorphic SNPs). The threshold chosen depends on the size of the study and the effect sizes expected. Power calculation software such as CaTS Power (Skol et al., 2006) or Quanto (Gauderman, 2002) can simplify power calculations for genetic association studies and inform the investigator of the allele frequency below which the study becomes severely underpowered. Minor allele frequency can be reported for each SNP using the `—freq` option in PLINK, and SNPs can be removed from the analysis using the `—maf` option, followed by a lower limit threshold. SNPs with frequency too low to yield reasonable statistical power (e.g., below 1%) may be removed from the analysis to lighten the computational and multiple testing correction burden. However, in studies with very large sample sizes, it may be beneficial to avoid removing these rare SNPs. Others have shown that nonsynonymous, possibly deleterious SNPs are on average rarer than synonymous SNPs that likely do not cause any adverse phenotypes (Gorlov et al., 2008).

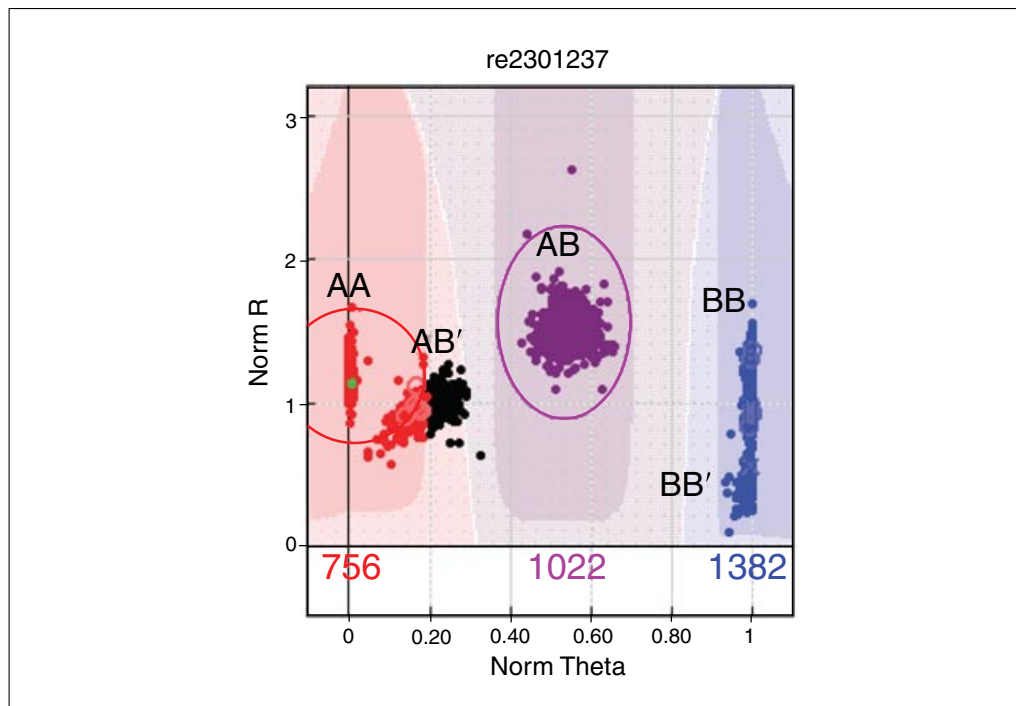


Figure 1.19.8 AB and BB individuals are split into subclusters AB and AB', BB and BB', while the AA cluster is unaffected. The AB/AB' split results in some AB samples miscalled as AA (diagnosed by Mendelian inconsistencies in the genotypes), as well as deviation from HWE due to excess homozygosity. Since only samples with at least one B allele demonstrate the splitting, one consistent explanation is the presence of a cryptic polymorphism near rs2301237 on a haplotype that contains the B allele. In this case, a second polymorphism (rs3114267) lies eight bases upstream from the typed polymorphism, and is in complete LD ($D' = 1$, $r^2 = 0.2$) with rs2301237.

Hardy-Weinberg equilibrium

Checking for Hardy-Weinberg Equilibrium (HWE) is one final step in the quality control analysis of markers in GWAS data. Under Hardy-Weinberg assumptions, allele and genotype frequencies can be estimated from one generation to the next. Departure from this equilibrium can be indicative of potential genotyping errors, population stratification, or even actual association to the trait under study (Wittke-Thompson et al., 2005). *UNIT 1.18* contains a very detailed description of the key principles and assumptions of HWE and how HWE is tested and applied in genetic association studies. HWE can be assessed using the `—hardy` option in PLINK. While departure from HWE can indicate potential genotyping error, disequilibrium can also result from a true association. It has been consistently noted that many more SNPs are out of HWE at any given significance threshold than would be expected by chance. SNPs severely out of HWE should therefore not be eliminated from the analysis,

but flagged for further analysis after the association analyses are performed.

Databases such as MySQL (see Table 1.19.1) can be very useful for joining association statistics with HWE statistics for easy reporting. It is also beneficial to examine HWE in controls separately, as disease-free controls should more closely follow the assumptions that lead to HWE than cases, and because some true associations are expected to be out of HWE. If multiple ethnicities are used in the same study, it is necessary to test for HWE within each group separately. SNPs that are highly associated with the trait of interest that also show highly significant departures from HWE, especially in controls, should be closely scrutinized. Typically HWE deviations toward an excess of heterozygotes reflect a technical problem in the assay, such as nonspecific amplification of the target region. On many GWAS platforms, the quantitative allelic signals at a marker, i.e., the intensity plot for the SNP,

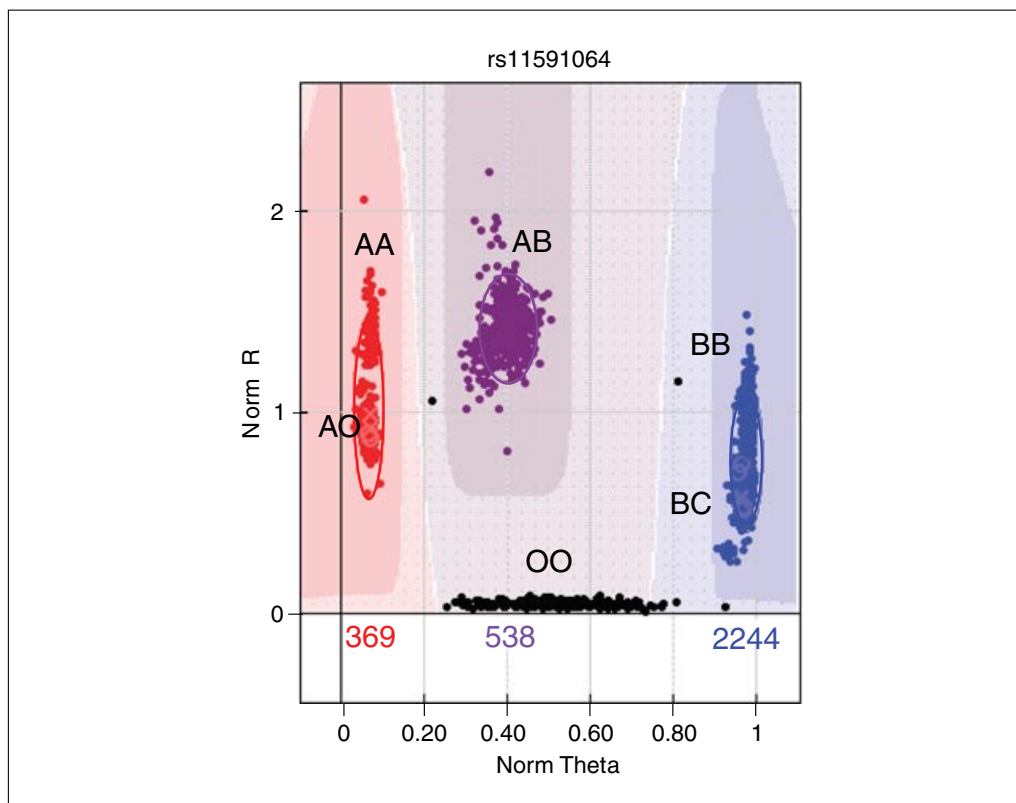


Figure 1.19.9 Unexpected number of clusters resulting in departure from HWE consistent with copy loss. Hemizygous individuals cluster at AO and BO. Individuals with homozygous deletions cluster at OO and their genotype calls are missing. The AB cluster remains intact, since these individuals are ipso facto diploid at the locus. Parent-parent-child Mendelian errors are present when at least one parent is hemizygous and produces hemizygous offspring. The deletion results in excess homozygosity. In this case, the “copy loss” appears to be a six-nucleotide insertion (rs71578153) coincident with rs11591064 that disrupts both A and B probes.

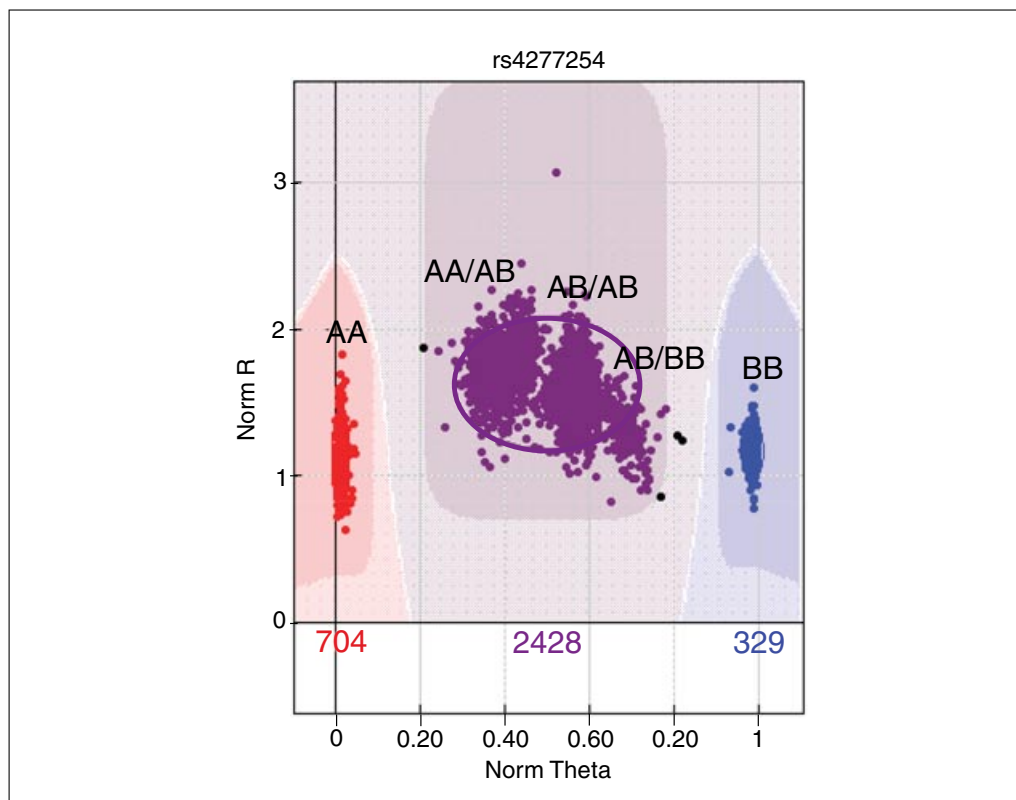


Figure 1.19.10 The five observed clusters are most consistent with a segmental duplication, although none is curated around the locus. A copy number variant would be expected to produce additional clusters above the AA and BB clusters (i.e., AAA and BBB), as opposed to the splits being confined to strictly the heterozygous clusters. Regardless, the artifact results in excess heterozygosity.

can be used to screen for a technical origin of the HWE deviation: null alleles can produce multimodal genotype clusters in the heterozygote clusters and one of the homozygote clusters (Fig. 1.19.8), or can produce an unexpected number of samples with no signal (Fig. 1.19.9); SNPs within CNVs or segmental duplications can produce clusters of genotypes intermediate between the three expected clusters of genotype (Fig. 1.19.10; Carlson et al., 2006). In the loci depicted in Figures 1.19.8 to 1.19.10, chi-square tests for HWE are rejected at p -values less than 10^{-80} , so these represent the most egregious examples of the aforementioned behavior. If no technical errors are detected, then a number of biologically plausible explanations exist for HWE deviations toward an excess of homozygotes: population stratification, assortative mating, and inbreeding, to name a few.

BATCH EFFECTS

Thousands of DNA samples are typically genotyped in a GWAS, which necessitates partitioning samples into small batches of sam-

ples processed in the lab together for genotyping (e.g., the set of samples on a 96-well plate). The precise size and composition of the sample batch depend on the array and lab process used. Systematic differences among the compositions of individuals in a batch (i.e., the case to control ratio or race/ethnicity of individuals on plates), and the within-plate accuracy and efficiency, can result in batch effects—apparent associations confounded by batch. The problem is in essence the same problem observed with population stratification—namely, that if there is an imbalance of cases and controls on a plate, and there are nonrandom (unknown) biases or inaccuracies in genotyping that differ from plate to plate, spurious associations will result.

Ideally, no batch effect will be present because individuals with different phenotypes, sex, race, and other confounders should be plated randomly, and because modern high-throughput genotyping technology is much more accurate, efficient, and consistent than earlier generations of GWAS assays. There are several approaches for examining a dataset for potential batch effects. One simple approach is

to calculate the average minor allele frequency and average genotyping call rate across all SNPs for each plate. Gross differences in either of these on any plates can easily be identified. Another method involves coding case/control status by plate followed by running the GWAS analysis testing each plate against all other plates. For example, the status of all samples on plate or batch 1 will be coded as case, while the status of every other sample is to be coded control. A GWAS analysis is to be performed (e.g., using the `—assoc` option in PLINK), and both the average *p*-value and the number of results significant at a certain threshold (e.g., $p < 1 \times 10^{-4}$) can be recorded. SNPs with low minor allele frequency (i.e., $< 5\%$) should be removed before this analysis is performed to improve the stability of test statistics. This procedure should be repeated for each plate or batch in the study. If any single plate has many more or many fewer significant results, or has an average *p*-value that deviates from 0.5 (under the null the average *p*-value will be 0.5 over many tests), then this batch should be further investigated for genotyping or composition problems. If batch effects are present, methods similar to those employed for population stratification (e.g., genomic control) may be used to mitigate the confounding effects.

EVALUATION OF QC AFTER ASSOCIATION ANALYSIS

After phenotypic association analysis, the quality control measures used should be evaluated. One method is to compare the observed number of statistically significant results (*p*-values) with the expected uniform distribution of *p*-values under the null hypothesis of no association. Too many significant results may indicate insufficient QC. Also, because no QC will catch all problematic SNPs, the intensity plots for statistically significant SNPs must be reviewed to make sure that there are no obvious clustering problems. Replication of results using different genotyping technology (such as TaqMan) and/or in another sample may be needed as well.

FUTURE DIRECTIONS

The QC pipeline developed by the eMERGE network has enabled a thorough analysis of the quality of the genome-wide genotype data generated on the ~17,000 samples. All of these data have been deposited in dbGaP along with corresponding quality con-

trol documents that describe all of the QC details for each dataset individually. Conducting QC in parallel at a coordinating center and study sites has been a tremendously valuable experience, as it led to a more thorough understanding since each group had to reconcile its results with others. Finally, this review has focused on the quality control of GWAS data, which typically only investigate common variants. Next-generation sequencing (NGS) of rare variants using short reads presents its own set of challenges in addition to those presented here. Future studies should consider issues unique to association studies using NGS.

ACKNOWLEDGEMENTS

This research was supported in part by NIH grants U01HG004608, U01HG004609, U01HG004610, U01HG04599, U01HG04603, U01HG004438, R01LM010040, and by the Intramural Research Program of the NIH, National Library of Medicine.

LITERATURE CITED

- 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.
- Barber, M.J., Mangravite, L.M., Hyde, C.L., Chasman, D.I., Smith, J.D., McCarthy, C.A., Li, X., Wilke, R.A., Rieder, M.J., Williams, P.T., Ridker, P.M., Chatterjee, A., Rotter, J.I., Nickerson, D.A., Stephens, M., and Krauss, R.M. 2010. Genome-wide association of lipid-lowering response to statins in combined study populations. *PLoS One* 5:e9763.
- Broman, K.W. 1999. Cleaning genotype data. *Genet. Epidemiol.* 17:S79-S83.
- Cardon, L.R. and Palmer, L.J. 2003. Population stratification and spurious allelic association. *Lancet* 361:598-604.
- Carlson, C.S., Smith, J.D., Stanaway, I.B., Rieder, M.J., and Nickerson, D.A. 2006. Direct detection of null alleles in SNP genotyping data. *Hum. Mol. Genet.* 15:1931-1937.
- Chanock, S.J., Manolio, T., Boehnke, M., Boerwinkle, E., Hunter, D.J., Thomas, G., Hirschhorn, J.N., Abecasis, G., Altshuler, D., Bailey-Wilson, J.E., Brooks, L.D., Cardon, L.R., Daly, M., Donnelly, P., Fraumeni, J.F. Jr., Freimer, N.B., Gerhard, D.S., Gunter, C., Guttmacher, A.E., Guyer, M.S., Harris, E.L., Hoh, J., Hoover, R., Kong, C.A., Merikangas, K.R., Morton, C.C., Palmer, L.J., Phimister, E.G., Rice, J.P., Roberts, J., Rotimi, C., Tucker, M.A., Vogan, K.J., Wacholder, S., Wijsman, E.M., Winn, D.M., and Collins, F.S. 2007. Replicating genotype-phenotype associations. *Nature* 447:655-660.

- 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.
- Daly, A.K., Donaldson, P.T., Bhatnagar, P., Shen, Y., Pe'er, I., Floratos, A., Daly, M.J., Goldstein, D.B., John, S., Nelson, M.R., Graham, J., Park, B.K., Dillon, J.F., Bernal, W., Cordell, H.J., Pirmohamed, M., Aithal, G.P., Day, C.P.; DILI-GEN Study; International SAE Consortium. 2009. HLA-B*5701 genotype is a major determinant of drug-induced liver injury due to flucloxacillin. *Nat Genet* 41:816-819.
- Devlin, B. and Roeder, K. 1999. Genomic control for association studies. *Biometrics* 55:997-1004.
- Devlin, B., Bacanu, S.A., and Roeder, K. 2004. Genomic Control to the extreme. *Nat. Genet.* 36:1129-1130.
- Dumitrescu, L.C., Ritchie, M.D., Brown-Gentry, K., Pulley, J.J., Basford, M., Denny, J., Oksenberg, J.R., Roden, D.M., Haines, J.L., and Crawford, D.C. 2010. Assessing the accuracy of observer-reported ancestry in a biorepository linked to electronic medical records. *Genet. Med.* In press.
- Frayling, T.M. 2007. Genome-wide association studies provide new insights into type 2 diabetes aetiology. *Nat. Rev. Genet.* 8:657-662.
- Gauderman, W.J. 2002. Sample size requirements for matched case-control studies of gene-environment interaction. *Stat. Med.* 21:35-50.
- Gorlov, I.P., Gorlova, O.Y., Sunyaev, S.R., Spitz, M.R., and Amos, C.I. 2008. Shifting paradigm of association studies: Value of rare single-nucleotide polymorphisms. *Am. J. Hum. Genet.* 82:100-112.
- Grady, B.J., Torstenson, E., Dudek, S.M., Giles, J., Sexton, D., and Ritchie, M.D. 2010. Finding unique filter sets in plato: A precursor to efficient interaction analysis in gwas data. *Pac. Symp. Biocomput.* 2010:315-326.
- Hindorf, L.A., Sethupathy, P., Junkins, H.A., Ramos, E.M., Mehta, J.P., Collins, F.S., and Manolio, T.A. 2009. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl. Acad. Sci. U.S.A.* 106:9362-9367.
- International HapMap consortium. 2003. The International HapMap Project. *Nature* 426:789-796.
- International HapMap Consortium. 2007. A second generation human haplotype map of over 3.1 million SNPs. *Nature* 449:851-861.
- Kathiresan, S., Willer, C.J., Peloso, G.M., Demissie, S., Musunuru, K., Schadt, E.E., Kaplan, L., Bennett, D., Li, Y., Tanaka, T., Voight, B.F., Bonnycastle, L.L., Jackson, A.U., Crawford, G., Surti, A., Guiducci, C., Burt, N.P., Parish, S., Clarke, R., Zelenika, D., Kubalanza, K.A., Morken, M.A., Scott, L.J., Stringham, H.M., Galan, P., Swift, A.J., Kuusisto, J., Bergman, R.N., Sundvall, J., Laakso, M., Ferrucci, L., Scheet, P., Sanna, S., Uda, M., Yang, Q., Lunetta, K.L., Dupuis, J., de Bakker, P.I., O'Donnell, C.J., Chambers, J.C., Kooner, J.S., Hercberg, S., Meneton, P., Lakatta, E.G., Scuteri, A., Schlessinger, D., Tuomilehto, J., Collins, F.S., Groop, L., Altshuler, D., Collins, R., Lathrop, G.M., Melander, O., Salomaa, V., Peltonen, L., Orho-Melander, M., Ordovas, J.M., Boehnke, M., Abecasis, G.R., Mohlke, K.L., and Cupples, L.A. 2009. Common variants at 30 loci contribute to polygenic dyslipidemia. *Nat. Genet.* 41:56-65.
- Klein, R.J., Zeiss, C., Chew, E.Y., Tsai, J.Y., Sackler, R.S., Haynes, C., Henning, A.K., Sangiovanni, J.P., Mane, S.M., Mayne, S.T., Bracken, M.B., Ferris, F.L., Ott, J., Barnstable, C., and Hoh, J. 2005. Complement factor H polymorphism in age-related macular degeneration. *Science* 308:385-389.
- Laurie, C., Mirel, D., Pugh, E., Bierut, L., Bhangale, T., Boehm, F., Caporaso, N., Edenburgh, H., Gabriel, S., Harris, E., Hu, F.B., Jacobs, K.B., Kraft, P., Landi, M.T., Lumley, T., Manolio, T.A., McHugh, C., Painter, J., Paschall, J., Rice, J.P., Rice, K.M., Zheng, X., Weir, B.S.; GENEVA Investigators. 2010. Quality control and quality assurance in genotypic data for genome-wide association studies. *Genet. Epidemiol.* 34:591-602.
- Link, E., Parish, S., Armitage, J., Bowman, L., Heath, S., Matsuda, F., Gut, I., Lathrop, M., and Collins, R. 2008. SLC6B1 variants and statin-induced myopathy: A genomewide study. 2008. *N. Engl. J. Med.* 359:789-799.
- Mailman, M.D., Feolo, M., Jin, Y., Kimura, M., Tryka, K., Bagoutdinov, R., Hao, L., Kiang, A., Paschall, J., Phan, L., Popova, N., Pretel, S., Ziyabari, L., Lee, M., Shao, Y., Wang, Z.Y., Sirotkin, K., Ward, M., Kholodov, M., Zbicz, K., Beck, J., Kimelman, M., Shevelev, S., Preuss, D., Yaschenko, E., Graeff, A., Ostell, J., and Sherry, S.T. 2007. The NCBI dbGaP database of genotypes and phenotypes. *Nat. Genet.* 39:1181-1186.
- Manolio, T.A. 2009. Collaborative genome-wide association studies of diverse diseases: Programs of the NHGRI's office of population genomics. *Pharmacogenomics* 10:235-241.
- Marchini, J., Cardon, L.R., Phillips, M.S., and Donnelly, P. 2004. The effects of human population structure on large genetic association studies. *Nat. Genet.* 36:512-517.
- McCarty, C., Chrisolm, R., Chute, C., Kullo, I., Jarvik, G., Larson, E., Li, R., Masys, D., Ritchie, M., Roden, D. et al. 2010. The eMERGE Network: A consortium of biorepositories linked to electronic medical records data for conducting genomic studies. *BMC Medical Genomics* In press.
- Miyagawa, T., Nishida, N., Ohashi, J., Kimura, R., Fujimoto, A., Kawashima, M., Koike, A., Sasaki, T., Tani, H., Otowa, T., Momose, Y., Nakahara, Y., Gotoh, J., Okazaki, Y., Tsuji, S., and Tokunaga, K. 2008. Appropriate data cleaning methods for genome-wide association study. *J. Hum. Genet.* 53:886-893.
- Newton-Cheh, C., Johnson, T., Gateva, V., Tobin, M.D., Bochud, M., Coin, L., Najjar, S.S., Zhao, J.H., Heath, S.C., Eyheramendy, S., Papadakis,

- K., Voight, B.F., Scott, L.J., Zhang, F., Farrall, M., Tanaka, T., Wallace, C., Chambers, J.C., Khaw, K.T., Nilsson, P., van der Harst, P., Polidoro, S., Grobbee, D.E., Onland-Moret, N.C., Bots, M.L., Wain, L.V., Elliott, K.S., Teumer, A., Luan, J., Lucas, G., Kuusisto, J., Burton, P.R., Hadley, D., McArdle, W.L.; Wellcome Trust Case Control Consortium, Brown, M., Dominiczak, A., Newhouse, S.J., Samani, N.J., Webster, J., Zeggini, E., Beckmann, J.S., Bergmann, S., Lim, N., Song, K., Vollenweider, P., Waeber, G., Waterworth, D.M., Yuan, X., Groop, L., Orho-Melander, M., Allione, A., Di Gregorio, A., Guarrera, S., Panico, S., Ricceri, F., Romanazzi, V., Sacerdote, C., Vineis, P., Barroso, I., Sandhu, M.S., Luben, R.N., Crawford, G.J., Jousilahti, P., Perola, M., Boehnke, M., Bonnycastle, L.L., Collins, F.S., Jackson, A.U., Mohlke, K.L., Stringham, H.M., Valle, T.T., Willer, C.J., Bergman, R.N., Morken, M.A., Döring, A., Gieger, C., Illig, T., Meitinger, T., Org, E., Pfeufer, A., Wichmann, H.E., Kathiresan, S., Marrugat, J., O'Donnell, C.J., Schwartz, S.M., Siscovick, D.S., Subirana, I., Freimer, N.B., Hartikainen, A.L., McCarthy, M.I., O'Reilly, P.F., Peltonen, L., Pouta, A., de Jong, P.E., Snieder, H., van Gilst, W.H., Clarke, R., Goel, A., Hamsten, A., Peden, J.F., Seedorf, U., Syvänen, A.C., Tognoni, G., Lakatta, E.G., Sanna, S., Scheet, P., Schlessinger, D., Scuteri, A., Dörr, M., Ernst, F., Felix, S.B., Homuth, G., Lorbeer, R., Reffellmann, T., Rettig, R., Völker, U., Galan, P., Gut, I.G., Herberg, S., Lathrop, G.M., Zelenika, D., Deloukas, P., Soranzo, N., Williams, F.M., Zhai, G., Salomaa, V., Laakso, M., Elosua, R., Forouhi, N.G., Völzke, H., Uitterwaal, C.S., van der Schouw, Y.T., Numans, M.E., Matullo, G., Navis, G., Berglund, G., Bingham, S.A., Kooner, J.S., Connell, J.M., Bandinelli, S., Ferrucci, L., Watkins, H., Spector, T.D., Tuomilehto, J., Altshuler, D., Strachan, D.P., Laan, M., Meneton, P., Wareham, N.J., Uda, M., Jarvelin, M.R., Mooser, V., Melander, O., Loos, R.J., Elliott, P., Abecasis, G.R., Caulfield, M., and Munroe, P.B. 2009. Genome-wide association study identifies eight loci associated with blood pressure. *Nat. Genet.* 41:666-676.
- Novembre, J., Johnson, T., Bryc, K., Kutalik, Z., Boyko, A.R., Auton, A., Indap, A., King, K.S., Bergmann, S., Nelson, M.R., Stephens, M., and Bustamante, C.D. 2008. Genes mirror geography within Europe. *Nature* 456:98-101.
- Patterson, N., Price, A.L., and Reich, D. 2006. Population structure and eigenanalysis. *PLoS Genet.* 2:e190.
- 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., Stephens, M., and Donnelly, P. 2000. Inference of population structure using multilocus genotype data. *Genetics* 155:945-959.
- Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A., Bender, D., Maller, J., Sklar, P., de Bakker, P.I., 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.
- Reich, D.E. and Goldstein, D.B. 2001. Detecting association in a case-control study while correcting for population stratification. *Genet. Epidemiol.* 20:4-16.
- Sherry, S.T., Ward, M.H., Kholodov, M., Baker, J., Phan, L., Smigielski, E.M., and Sirotkin, K. 2001. dbSNP: The NCBI database of genetic variation. *Nucleic Acids Res.* 29:308-311.
- Simon-Sanchez, J., Scholz, S., Fung, H.C., Matarin, M., Hernandez, D., Gibbs, J.R., Britton, A., de Vrieze, F.W., Peckham, E., Gwinn-Hardy, K., Crawley, A., Keen, J.C., Nash, J., Borgaonkar, D., Hardy, J., and Singleton, A. 2007. Genome-wide SNP assay reveals structural genomic variation, extended homozygosity and cell-line induced alterations in normal individuals. *Hum. Mol. Genet.* 16:1-14.
- Skol, A.D., Scott, L.J., Abecasis, G.R., and Boehnke, M. 2006. Joint analysis is more efficient than replication-based analysis for two-stage genome-wide association studies. *Nat. Genet.* 38:209-213.
- Tang, H., Quertermous, T., Rodriguez, B., Kardia, S.L., Zhu, X., Brown, A., Pankow, J.S., Province, M.A., Hunt, S.C., Boerwinkle, E., Schork, N.J., and Risch, N.J. 2005. Genetic structure, self-identified race/ethnicity, and confounding in case-control association studies. *Am. J. Hum. Genet.* 76:268-275.
- Thompson, J.F., Hyde, C.L., Wood, L.S., Paciga, S.A., Hinds, D.A., Cox, D.R., Hovingh, G.K., and Kastelein, J.J. 2009. Comprehensive whole-genome and candidate gene analysis for response to statin therapy in the Treating to New Targets (TNT) cohort. *Circ. Cardiovasc. Genet.* 2:173-181.
- Willer, C.J., Sanna, S., Jackson, A.U., Scuteri, A., Bonnycastle, L.L., Clarke, R., Heath, S.C., Timpson, N.J., Najjar, S.S., Stringham, H.M., Strait, J., Duren, W.L., Maschio, A., Busonero, F., Mulas, A., Albai, G., Swift, A.J., Morken, M.A., Narisu, N., Bennett, D., Parish, S., Shen, H., Galan, P., Meneton, P., Herberg, S., Zelenika, D., Chen, W.M., Li, Y., Scott, L.J., Scheet, P.A., Sundvall, J., Watanabe, R.M., Nagaraja, R., Ebrahim, S., Lawlor, D.A., Ben-Shlomo, Y., Davey-Smith, G., Shuldiner, A.R., Collins, R., Bergman, R.N., Uda, M., Tuomilehto, J., Cao, A., Collins, F.S., Lakatta, E., Lathrop, G.M., Boehnke, M., Schlessinger, D., Mohlke, K.L., and Abecasis, G.R. 2008. Newly identified loci that influence lipid concentrations and risk of coronary artery disease. *Nat. Genet.* 40:161-169.
- Wittke-Thompson, J.K., Pluzhnikov, A., and Cox, N.J. 2005. Rational inferences about departures from Hardy-Weinberg equilibrium. *Am. J. Hum. Genet.* 76:967-986.
- Zhang, F., Wang, Y., and Deng, H.W. 2008. Comparison of population-based association study methods correcting for population stratification. *PLoS One* 3:e3392.

INTERNET RESOURCES

<http://censtats.census.gov/data/WI/1605549675.pdf>
Census 2000. Profile of Demographic Characteristics, Marshfield, Wisconsin.

<http://pngu.mgh.harvard.edu/~purcell/plink/>
Illumina Technical Note: "TOP/BOT" Strand and "A/B" Allele (2009).

http://www.illumina.com/Documents/products/technotes/technote_gen_call_data_analysis_software.pdf
Illumina GenCall Data Analysis Software (2008).

<http://www.R-project.org>

R Development Core Team: R: A language and environment for statistical computing. ISBN 3900051070, Vienna, Austria: R Foundation for Statistical Computing (2005).

<http://pritch.bsd.uchicago.edu/structure.html>.
STRUCTURE (2009).

https://www.mc.vanderbilt.edu/victr/dcc/projects/acc/index.php/Visualizing_relatedness
Turner, S.D. 2009. Visualizing sample relatedness in a GWAS using PLINK and R.