



NIH PUBLIC ACCESS

Author Manuscript

*Genet Epidemiol.* Author manuscript; available in PMC 2014 July 07.

Published in final edited form as:

*Genet Epidemiol.* 2013 November ; 37(7): 666–674. doi:10.1002/gepi.21747.

## The Value of Statistical or Bioinformatics Annotation for Rare Variant Association with Quantitative Trait

Andrea E. Byrnes<sup>1</sup>, Michael C. Wu<sup>1</sup>, Fred A. Wright<sup>1</sup>, Mingyao Li<sup>2</sup>, and Yun Li<sup>1,3,4</sup><sup>1</sup>Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina 27599<sup>2</sup>Department of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA 19104<sup>3</sup>Department of Genetics, University of North Carolina, Chapel Hill, North Carolina 27599<sup>4</sup>Department of Computer Science, University of North Carolina, Chapel Hill, North Carolina 27599

### Abstract

In the past few years, a plethora of methods for rare variant association with phenotype have been proposed. These methods aggregate information from multiple rare variants across genomic region(s), but there is little consensus as to which method is most effective. The weighting scheme adopted when aggregating information across variants is one of the primary determinants of effectiveness. Here we present a systematic evaluation of multiple weighting schemes through a series of simulations intended to mimic large sequencing studies of a quantitative trait. We evaluate existing phenotype-independent and -dependent methods, as well as weights estimated by penalized regression approaches including Lasso, Elastic Net and SCAD. We find that the difference in power between phenotype-dependent schemes is negligible when high quality functional annotations are available. When functional annotations are unavailable or incomplete, all methods suffer from power loss; however, the variable selection methods outperform the others at the cost of increased computational time. Therefore, in the absence of good annotation, we recommend variable selection methods (which can be viewed as “statistical annotation”) on top regions implicated by a phenotype independent weighting scheme. Further, once a region is implicated, variable selection can help to identify potential causal SNPs for biological validation. These findings are supported by an analysis of a high coverage targeted sequencing study of 1898 individuals.

### Keywords

rare variants; association; weighting; variable selection; variant annotation

---

Correspondence Should Be Addressed To: Yun Li. Mailing Address: Department of Genetics, Campus Box 7264, University of North Carolina, Chapel Hill, North Carolina 27599, Phone: (919) 843 2832, Fax: (919) 843 4682, [yunli@med.unc.edu](mailto:yunli@med.unc.edu).

**Supplementary Data:** Supplemental Data which include four Supplementary Figures can be found with the article.

## Introduction

Recent studies have shown that rare variants may be important to the underlying etiology of complex traits [Cohen et al., 2004; Dickson et al., 2010; Gorlov et al., 2008; Haase et al., 2012; Nelson et al., 2012; Zawistowski et al., 2010] and that they may account for part of the “missing heritability” [Eichler et al., 2010; Gibson 2010; Maher 2008; Manolio et al. 2009] left by Genome-wide Association Studies (GWAS). Conventional association analysis methods, which evaluate each variant independently of all others, lack the statistical power to evaluate rare variants given the sample size of sequencing data currently available. However, there is increasing evidence that the combined effects of rare variants in the same exon, gene, region or biological pathway can be used to elucidate complex phenotypes [Cohen et al., 2004; Nejentsev et al., 2009; Sanna et al., 2011]. Where the effect size of a single variant may not be large enough to detect with the sample sizes available, a collection of variants with small effect size, taken together, may be detectable. In order to explore the potential effects of rare variants in present-day genomic data, a large number of methods [Bacanu et al., 2011; Cheung et al., 2012; Lee et al., 2012; Li and Leal 2008; Li et al., 2010a; Madsen and Browning 2009; Mao et al., 2012; Neale et al., 2011; Price et al., 2010; Tzeng et al., 2011; Wu et al., 2011; Xu et al., 2012; Yi et al., 2011] for aggregating information across variants have emerged. However there is little consensus on which method is most effective. The weighting scheme adopted when aggregating across variants is an important consideration, as is the use of functional or bioinformatics information when available.

We present an evaluation of multiple weighting schemes through a series of simulations. We evaluate several existing phenotype-independent [Cohen et al., 2004; Madsen and Browning 2009; Morgenthaler and Thilly 2007] and -dependent weighting schemes [Wu et al., 2011; Xu et al., 2012], as well as weighting schemes determined by linear regression, penalized regression and variable selection methods, including Lasso [Tibshirani 1996], Elastic Net [Zou and Hastie 2005] and SCAD [Xie and Huang 2009]. We conduct simulations under a variety of scenarios with different numbers of true causal variants, mixtures of direction of effect and availability of functional information, mimicking sequencing studies of a quantitative trait. We then apply each of these methods to a set of high coverage targeted sequencing data [Nelson et al., 2012] of 1898 individuals from the CoLaus population-based cohort [Firmann et al., 2008].

## Materials and Methods

Over the last few years, numerous sensible weighting schemes have been proposed. In most of these methods a genomic region or variant set is assigned a weighted sum over the variants meant to describe the burden of potentially influential variants carried by each individual. We call this weighted sum  $S_i$ . Further, we assume there are  $N$  individuals under study, indexed by  $i$ , and for each individual we have  $M$  variants in the region or variant set, indexed by  $j$ .

## Phenotype-Independent Weighting Schemes

First, we examine three approaches that are independent of the observed phenotype. The first of these is a simple indicator of whether or not rare variants (minor allele frequency,  $MAF < 0.01$ ) are present in the region [Cohen et al., 2004]. That is,

$$S_i = I \left( \sum_{j=1}^M I(\hat{q}_j < Q) x_{ij} > 0 \right)$$

where  $x_{ij}$  is the number of minor alleles observed for individual  $i$  at variant  $j$ .  $\hat{q}_j = \frac{\sum_{i=1}^N x_{ij} + 1}{2N + 2}$  is the estimated MAF of variant  $j$  in the data with pseudo counts and  $Q$  is the MAF threshold. In this work, we consider  $Q = 0.05$ .

Second, we examine a count approach which assigns a higher score to individuals carrying a larger number of rare alleles [Morgenthaler and Thilly 2007];

$$S_i = \sum_{j=1}^M I(\hat{q}_j < Q) x_{ij}.$$

with  $x_{ij}$  being the count of rare alleles for individual  $i$  at variant  $j$  and  $\hat{q}_j$  being the estimated MAF, as defined above.

We also consider the approach proposed by Madsen and Browning [Madsen and Browning 2009] where the weight for variant  $j$  is a function of the minor allele frequency (MAF):

$$S_i = \sum_{j=1}^M \xi_j x_{ij}, \text{ where } \xi_j = \frac{1}{\sqrt{N \times \hat{q}_j \times (1 - \hat{q}_j)}}$$

with  $x_{ij}$  and  $\hat{q}_j$  as above. In the original Madsen and Browning framework for case-control studies, MAFs are estimated using controls only. However, in this paper, the outcome of interest is quantitative and we estimate MAF using the entire sample, which makes the method phenotype-independent in this context.

## Phenotype-Dependent Weighting Schemes

We also consider phenotype-dependent regression-based methods. First, we examine the performance of marginal regression coefficients. That is, we fit the simple linear regression model  $Y = x_j \beta_j + \varepsilon$  for each variant  $j$  separately and independently and then take the fitted values  $\hat{\beta}_j$  to be our weights.

$$S_i = \sum_{j=1}^M \xi_j x_{ij}, \text{ where } \xi_j = \tilde{\beta}_j, \text{ the MLE of } \beta \text{ for the model above.}$$

Though imperfect, this weighting scheme allows investigators to test for associations with multiple rare variants in cases where  $N < M$  and begin to follow up on individual variants that may potentially be of interest.

Second, we consider weights from ordinary multiple regression, modeling all of the  $M$  variants simultaneously. That is, we fit the model  $Y = X\beta + \varepsilon$ , where the  $(i, j)^{\text{th}}$  element of the matrix  $X = x_{ij}$ , the minor allele count for individual  $i$  at variant  $j$ . We then take  $S_i$  to be as above, with the fitted values from this multiple regression,  $\beta_i = \tilde{\xi}_j$  [Lin and Tang 2011; Xu et al., 2012].

We also consider weights from several variable selection methods. Such methods are appealing since we expect the majority of rare variants not to influence the quantitative trait of interest. Use of penalized regression is therefore expected to reduce the number of non-zero weights. Similar strategies were recently proposed in the context of rare variant association testing [Turkmen & Lin, 2012; Zhou, Sehl, Sinsheimer, & Lange, 2010]. In penalized regression, we solve for the  $\tilde{\beta}$ 's which best fit the data, subject to some constraint(s) or penalty. That is, instead of minimizing the sum of squared error,  $(Y - \beta X)'(Y - \beta X)$ , we aim to minimize the sum of squared errors and an additional penalty term,  $(Y - \beta X)'(Y - \beta X) + P(\lambda, \beta)$ . In general, the greater the number of parameters included in the model, the greater the penalty. A number of penalty functions have been proposed and extensively studied in the recent statistical literature [Heckman and Ramsay 2000; Hesterberg et al., 2008; Kyung et al., 2010; Wu and Lange 2008]. Of these, we chose three: the Lasso which imposes a linear penalty [Tibshirani 1996], Elastic Net (EN) which imposes a quadratic penalty [Zou and Hastie 2005] and SCAD which is designed to penalize smaller coefficients more heavily than larger coefficients [Xie and Huang 2009].

For Lasso and SCAD, only one tuning parameter,  $\lambda$ , is required. We used the R packages *lars* [Efron, Hastie, Johnstone, & Tibshirani, 2004] and *ncvreg* [Breheny & Huang, 2011] with default parameter values, which is to choose the optimal  $\lambda$  among a grid of 100 possible values equally spaced on the log-scale. For Elastic Net, there are two tuning parameters, one for the linear component and one for the quadratic component. The linear term,  $\lambda_1$ , is chosen in the same way as the  $\lambda$  parameter for the Lasso and SCAD methods, discussed above. The quadratic parameter,  $\lambda_2$ , was set to 1 in all simulations and for the real data. We used the R package *elasticnet* to fit the EN models [Zou & Hastie, 2005]. After model fitting, we then use estimated coefficients from each of these variable selection methods as weights. The number of non-zero coefficients included is upper-bounded by 100 for each of these schemes throughout this work.

Under each weighting scheme examined, we determine the significance of a genomic region using a score test of the following form:

$$U = \sum_{i=1}^N (Y_i - \bar{Y}) S_i \text{ where}$$

$S_i = \sum_{j=1}^M \xi_j x_{ij}$  in which  $N$  is the number of individuals under study, and  $Y_i$  is the quantitative trait value for the  $i^{\text{th}}$  individual.  $S_i$  is the genetic score for the  $i^{\text{th}}$  individual, a weighted sum across multiple variants. Specifically,  $x_{ij}$  is the number of minor alleles observed for individual  $i$  at variant  $j$  where  $x_{ij}$  are not normalized.  $M$  is the number of variants in the region under study (discovered through sequencing in our context) and  $\xi_j$  is the weight of variant  $j$  under one of the above weighting schemes. The analytical distribution for this statistic is not generally known in this context, so significance must be assessed empirically by permutation.

Additionally, we apply the similarity-based method SKAT [Wu et al., 2011] to each of our simulated data sets and the real data set for comparison. We use weights based on the default Beta distribution implemented in the SKAT package, version 0.79. We will comment in the Discussion section on the conceptual differences between the weighting schemes we consider in this work and the SKAT methodology.

### Simulation Setup

We simulate 45,000 chromosomes for a series of 100 50Kb regions with a coalescent model [Schaffner et al. 2005] that mimics linkage disequilibrium (LD) in real data, accounts for variations in local recombination rates and models population history consistent with the CEU samples. We then randomly select 2,000 simulated chromosomes (forming 1,000 diploid individuals) to mimic a large sequencing study. For each region, we simulate one single pool of 45,000 chromosomes instead of multiple pools of 2,000 chromosomes so that the causal variants in each region can be determined by population MAFs (MAFs calculated using the entire population of 45,000 chromosomes) and thus retained across replicates from the same region. We assume only rare variants ( $0.001 < \text{population MAF} < 0.05$ ) influence the value of the quantitative trait and we randomly select  $m$  variants that truly influence the quantitative trait value. For each variant, we independently assign the direction of influence according to  $r$ , the probability that a causal variant will increase the trait value. Following Wu et al [Wu et al., 2011], we then simulate quantitative traits under the null model:

$$y_i = 0.5E_{1i} + 0.5E_{2i} + \varepsilon_i \quad (\text{null model})$$

where  $E_{1i}$ ,  $E_{2i}$  and  $\varepsilon_i$  are independent with  $E_{1i} \sim \text{Bernoulli}(0.5)$  to mimic a binary covariate,  $E_{2i} \sim \text{Normal}(0,1)$  to mimic a continuous covariate, and  $\varepsilon_i \sim \text{Normal}(0,1)$ . We also simulate quantitative traits under an alternative model:

$$y_i = 0.5E_{1i} + 0.5E_{2i} + \sum_{j=1}^m \beta x_{ij}^C + \varepsilon_i \quad (\text{Alternative Model})$$

where  $\beta_j = r_j |k \times F(MAF_j)|$  and  $r_j = 1$  with probability  $r$  and  $r_j = -1$  with probability  $(1-r)$ .

$E_{1i}$ ,  $E_{2i}$  and  $\varepsilon_i$  are as before,  $j$  indexes the truly causal variants and  $x_{ij}^C$  is the number of minor alleles individual  $i$  has at causal variant  $j$ . The link function  $F$  takes one of the following forms:

$$F_{\log}(q) = k \times \log(q), F_{\text{logit}}(q) = k \times \log\left(\frac{q}{1-q}\right), F_{MB}(q) = k \times \frac{1}{\sqrt{q(1-q)}}$$

where  $N$  is the number of individuals sequenced. We call the first link function log, the second logit, and the third Madsen-Browning (MB). In addition, we also consider  $F_{\text{random}}(q)$ , a random value chosen from the *exponential*(1) distribution, independent of  $q$  and multiplied by  $k$ . The constant  $k$  is a scaling factor to control the magnitude of the change in quantitative trait due to truly causal genetic variants. In our simulations  $k$  is set to 0.2, which keeps the heritability  $h^2$ , between 0.1% and 2.5%. Complex human quantitative traits are thought to have heritability estimates in this range [Manolio et al., 2009]. In the Results section, we report the results for the logit link function; results for all four link functions are given in the Supplementary materials.

To assess significance in each simulated setting, score test statistic from each weighting scheme is compared to the empirical distribution of the test statistic obtained under the null simulations. We assess the significance of each test at the  $\alpha=0.01$  level using the empirical null distribution, which we approximate using 100,000 data sets simulated under the null hypothesis of no variant contributing to the quantitative trait.

**Simulation of Data Sets under the Null Hypothesis**—For each of the 100 regions we simulate, we randomly select 100 samples of 2,000 chromosomes (forming 1,000 diploid individuals). We then assign quantitative trait values under the null model specified above. Using these  $100 \times 100=10,000$  data sets simulated under the null hypothesis, we obtain the empirical null distribution of the test statistics for each method.

**Simulation of Data Sets under Different Alternative Hypotheses**—For each choice of  $r$ ,  $m$  and  $F(\cdot)$ , we select 2,000 chromosomes from the population of 45,000 chromosomes again via simple random sampling. Again, we randomly pair these chromosomes to form diploid individuals and replicate 100 times for each region. For each replicate, we randomly select  $m$  rare variants to be causal. Each causal variant is assigned a direction in which to exert its effect (positive with probability  $r$  and negative with probability  $1-r$ ).

**Simulation of “Good” Functional Annotation**—In each simulated data set, we annotate variants as “functional” or “non-functional”. We assume that we have a reasonably

good bioinformatics tool such that a true causal variant has 90% probability to be annotated as “functional”. Even a perfect bioinformatics tool can only predict functionality, not causality or association with a particular trait of interest. Because of this, we annotate an additional random number of  $W$  non-causal variants as “functional”. Kryukov and colleagues [Kryukov et al., 2009] have estimated that approximately one third of *de novo* missense mutations (that would be predicted as functional by a sensible bioinformatics tool) have no effect on phenotypic traits. We therefore used 1/3 as the lower bound for the fraction of non-causal variants annotated and simulated  $W \sim N(25,5)$ , rounded to the nearest integer. We evaluate the performance of each of these weighting schemes both using all variants without the help of the bioinformatics tool, and using only the “functional” variants annotated. Under the null distribution,  $W$  variants are selected at random.

**Simulation of GWAS Data Sets**—We use the same choice of causal variants in each region as in the simulated sequencing data. Consequently, the direction of association and true effect size of each of these are unchanged. In order to simulate GWAS SNPs, we select 1000 chromosomes from the total 45,000 to mimic the 1000 Genomes [Abecasis et al., 2012] sample. The simulated 1000 Genomes sample is used to define LD, based on which GWAS SNPs are selected. For each region, we choose 75 GWAS SNPs consisting of the first 70 tagSNPs (SNPs with the highest number of LD buddies where an LD buddy is a SNP with which the  $r^2 > 0.8$ ) and 5 SNPs at random from the remaining set of SNPs, mimicking the Illumina Omni5 or Affymetrix Axiom high-density SNP genotyping platforms.

## Results

### In the Absence of a Bioinformatics Tool

Throughout our simulations, we observe several consistent patterns. First, when we apply these methods in the absence of a Bioinformatics tool (thus, all variants are included in analysis), variable selection schemes (most noticeably Lasso and EN) outperform other methods, including SKAT, in nearly all situations (notable exceptions are discussed below). For example, under the simulated setting of 10 causal variants, among which we expect to five increase quantitative trait value, the power is 80.0% and 83.7% for Lasso and EN, and is 0.4%, 7.3%, 7.6%, 43.2%, 25.3%, 60.5%, 41.3%, and 46.6% for Indicator, Count, Madsen-Browning, Marginal Regression, Multiple Regression, SCAD, SKAT (all variants), and SKAT (rare variants only) respectively (Figure 1a). Under the simulated setting of 50 causal variants among which 40 are expected to increase quantitative trait value, power is 100% for both Lasso and EN, and is 0.03%, 0.19%, 0.07%, 99.63%, 100%, 100%, 96.9%, and 98.5% for Indicator, Count, Madsen-Browning, Marginal Regression, Multiple Regression, SCAD, SKAT (all variants), and SKAT (rare variants only) respectively (Figure 1b).

### In the Presence of a Good Bioinformatics Tool

In the presence of a good bioinformatics tool (as introduced in the Methods section) the power increases for each of the methods previously discussed. Most notably, the phenotype-independent methods show a substantial gain in power once the bioinformatics tool is

applied. For example, under the simulated setting of 10 causal variants, among which five are expected to increase the quantitative trait value, the power is 99.83% and 99.80% for Lasso and EN, and is 23.91%, 17.15%, 18.85%, 97.87%, 99.73%, 99.76%, 98.49%, and 98.34% for Indicator, Count, Madsen-Browning, Marginal Regression, Multiple Regression, SCAD, SKAT (all variants), and SKAT (rare variants only) respectively (Figure 2a). Under the simulated setting of 50 causal variants, among which 40 increase quantitative trait value, power is 100% for both Lasso and EN, and is 99.38%, 98.89%, 96.51%, 100%, 100%, 100%, 100%, and 100% for Indicator, Count, Madsen-Browning, Marginal Regression, Multiple Regression, SCAD, SKAT (all variants), and SKAT (rare variants only) respectively (Figure 2b). Although power increases for all methods, the relative performance of the methods changes little from that under the absence of a bioinformatics tool.

### Effect of $m$ (the Number of Causal Variants) and $r$ (% of Positive Causal Variants)

As the number of true causal variants ( $m$ ) increases, so does power for all methods. This is to be expected since adding more causal variants increases the signal-to-noise ratio. When the number of true causal variants is very small, none of the methods have adequate power. Interestingly, it is in these situations where  $m$  is very small that SKAT manifests its advantage over other methods examined. As  $r$  gets smaller (that is, the probability that a causal variant will contribute positively to the quantitative trait values gets smaller), the power of the phenotype-independent methods decreases. For example, the phenotype-independent methods have close to 0 power when  $r=0.05$ ; while the phenotype-dependent methods are relatively unaffected by changing values of  $r$  (Figure 1a and Figure 2a). We also observe a slight dip in power in all of the phenotype-dependent schemes when  $r=0.5$  and no bioinformatics information is used (Figure 1a), which is to be expected since the signals from different directions are canceling one another. Similar trends are seen in all simulations with all four link functions (shown in supplementary materials).

### Weight Estimation Accuracy for Individual Variants

Table 1 shows the correlation between the true and estimated values of the weights for each method under the simulation settings in which the number of truly causal variants,  $m$ , is 10 and the proportion of variants contributing in the positive direction,  $r$ , is 80%. Of note, the correlation between true and estimated weights increases for all methods with the addition of bioinformatics filtering. The Elastic Net and Lasso yield the highest correlations between estimated and true weights, both in situations where we restrict to variants that are likely to be functional (Pearson correlations of 0.285 and 0.355), and when we do not (Pearson correlations of 0.744 and 0.778).

### Identification of Individual Causal Variants

When using variable selection schemes, we have the opportunity to identify individual causal variants within the region or variant set under study. Figure 3 illustrates the accuracy with which the causal variant(s) can be identified by each weighting scheme. Note that the causal variant(s) are not always 100% identified, but in many cases, the causal variant, or a variant in high LD ( $r^2 > 0.8$ ), have estimated non-zero weights. For example, if we fix  $m=10$ ,  $r=0.8$  and the logit link function, without considering LD buddies, we need to



consider the top 696 (109 and 12) variants in order to detect 90% (60%, 30%) of the causal variants using EN (Figure 3a); taking LD buddies into consideration, the numbers decrease to 378 (14 and 4) (Figure 3b). When we also consider functional information we consider fewer variants and narrow the field to include a higher proportion of truly causal variants. In this case, we need to consider the top 408 (16 and 4) variants in order to detect 90% (60%, 30%) of the causal variants (Figure 3c) without considering LD buddies; with LD buddies taken into consideration, the numbers decrease to 374 (13 and 3) (Figure 3d).

### Results with GWAS Data Sets

Studies that sequence a portion or the entirety of the genome are becoming increasingly common, but still much more GWAS data exist than sequencing data. Imputation has been shown to accurately predict genotypes at untyped variants from GWAS data in a variety of circumstances [Auer et al., 2012; de Bakker et al., 2008; Li et al., 2010a; Li et al., 2009b; Liu et al., 2012; Marchini and Howie 2010]. Using our simulated GWAS data and simulated reference, we observe that variable selection can improve power for GWAS data as well. However, the power is consistently lower than that under the sequencing setting due to the imperfect rescue of information through imputation (comparing Figure 1 with Supplementary Figure 3). In our simulations, the imputation accuracy is 99.66% for all variants and 99.98% for rare variants, but most of the inaccuracies are due to missed rare variants. In fact, among variants with  $MAF < 0.001$  nearly all inaccuracies are due failure to identify the minor allele. Specifically, the squared Pearson correlation between the imputed genotypes (continuous, ranging from 0 to 2) and the true underlying genotypes (coded as 0, 1 and 2) is only 0.2397 for variants with  $MAF < 0.001$ . Supplementary Figure 3 shows the relative power of these weighting schemes over a range of  $r$  (Supplementary Figure 3a) and  $m$  (Supplementary Figure 3b).

### Results with Real Data Set

Of the over 6,000 individuals in the CoLaus cohort [Firmann et al., 2008], 1,898 had recorded total cholesterol and targeted sequence data in 202 drug target genes [Nelson et al., 2012]. Sequencing was done at moderately high coverage (with median coverage 27X) and genotype calls were obtained using *SOAP-SNP* [Li et al., 2009a]. Sporadic missing genotypes were imputed with *MaCH* [Li et al., 2010b]. One gene previously known to be associated with total cholesterol in these data is used as a positive control. We test each of the 172 autosomal genes with and without removing nonfunctional variants using *ANNOVAR* [Wang et al., 2010]. For each method, we estimate weights in association with total cholesterol and, for the methods that accommodate covariates, we adjust for age, age<sup>2</sup>, sex and the first five principal components. For the phenotype-independent methods, no covariate adjustment is performed and significance is assessed by permutation of the  $Y_i$ 's. For methods allowing covariates (marginal and multiple regression, Lasso, EN and SCAD), permutation of outcomes alone is not appropriate. For these methods, we fit a regression model,  $Y_i \sim Z_i$ , where  $Z$  is the matrix of covariates and then obtain residuals,  $\varepsilon_i$ . The  $\varepsilon_i$ 's are then randomly permuted to obtain a set of  $\varepsilon_i^*$ 's, the permuted residuals. For each permutation, we fit the model  $\varepsilon_i^* \sim X_i$  in order to re-estimate the weights  $\xi_j$  and scores  $S_i$  as in [Davidson and Hinkley 1997]. We do 10,000 such permutations and, from these, obtain a null distribution of statistics with which to assess significance. Since SKAT produces

analytical p-values shown to preserve type I error [Wu et al., 2011], we use the SKAT analytical p-values without permutation.

When all variants regardless of bioinformatics prediction are included, the variable selection methods Lasso and EN yield the smallest p-values compared to other methods for the previously implicated gene. However, the previously implicated gene is not the most significant among the 172 genes tested. Using *ANNOVAR* annotations [Wang et al., 2010], we restrict to non-synonymous variants in coding regions of the genome only. When considering only these functional variants, most weighting schemes identify the correct gene with highly significant p-values (Table 2 and Supplementary Figure 4).

## Discussion

In summary, through extensive simulation studies with varying number, model, and direction of causal variant(s) contributing to a quantitative trait, we find that functional annotations derived from good set of bioinformatics tools can substantially boost power for rare variant association testing. In the absence of good bioinformatics tools, “statistical” annotation based on phenotype-dependent weighting of the variants, particularly through variable selection based methods to both select potentially causal/associated variants and estimate their effect sizes, manifests advantages. This observation holds for both sequencing-based studies or studies based on a combination of genotyping, sequencing, and imputation. We also find supporting evidence from application to a real sequencing-based data set.

The price one has to pay for adopting phenotype-dependent methods is the necessity of permutation, which can be easily performed through permuting of residuals for the analysis of quantitative traits [Davidson and Hinkley 1997; Lin 2005] or using the BiasedUrn method [Epstein et al., 2012] recently proposed for binary traits. This, in turn, increases computational costs. Therefore, we recommend primarily using phenotype-dependent weighting for refining the level of significance. That is, we recommend applying phenotype-dependent weighting only to genomic regions or variant sets that have strong evidence of association (but not necessarily reaching genome-wide significance) from methods that do not require permutation (for example, SKAT [Wu et al., 2011]).

We note that testing over a region by aggregating information across variants is a different task from estimating effect sizes of individual variant (as measured by the variant weights in our work). Perfection in the latter (that is, being able to estimate weights for each individual variants accurately) leads to perfection in the former (that is, maximal testing power over the region harboring those variants); but not vice versa. Based on our simulations where we know the true contribution (effect size) of each individual variant, we find that individual effect sizes cannot be well estimated (Pearson correlation between true and estimated effect sizes  $< 0.5$  even for the best variable selection based methods). However, these methods can still increase power of region or variant set association analysis without accurate estimation of individual variant effect sizes. In addition, these methods are able to identify the vast majority of the causal variants, particularly when LD buddies are considered.

In this paper, we mainly consider aggregation of information at the genotype level (where we first obtain a regional genotype score via a weighted sum of genotype scores for individual variants and then assess the association between the regional genotype score and the phenotype of interest), which underlies the largest number of rare variant association methods published. In contrast, there are methods that aggregate information at the effect size level (for example, SKAT [Wu et al., 2011] where the final regional score test statistic is a weighted sum of the test statistics for individual variants) or at the p-value level, for example in [Cheung et al., 2012]. Our comparisons with SKAT suggest that the same conclusions apply to aggregation methods at levels other than genotype.

Lastly, although one could potentially argue that the phenotype-dependent methods require an undesirable computing-power trade-off in the presence of good bioinformatics tools, in practice, we rarely (if ever) get perfect bioinformatics tools. In addition, even perfect bioinformatics tools can only predict functionality but NOT causality or association with particular phenotypic trait(s) of interest. Therefore, we view that the application of “statistical annotation” through phenotype-dependent weighting, particularly using variable selection based methods, to top regions or variant sets implicated by computationally efficient phenotype-independent methods, is valuable.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

We thank GlaxoSmithKline, especially Drs. Margaret G. Ehm, Matthew R. Nelson, Li Li, and Liling Warren for sharing the targeted sequencing data. We also thank our CoLaus collaborators for providing the phenotypic data. The research is supported by R01HG006292, R01HG006703 (awarded to Y.L.), and R01HG004517, R01HG005854 (to M.L.).

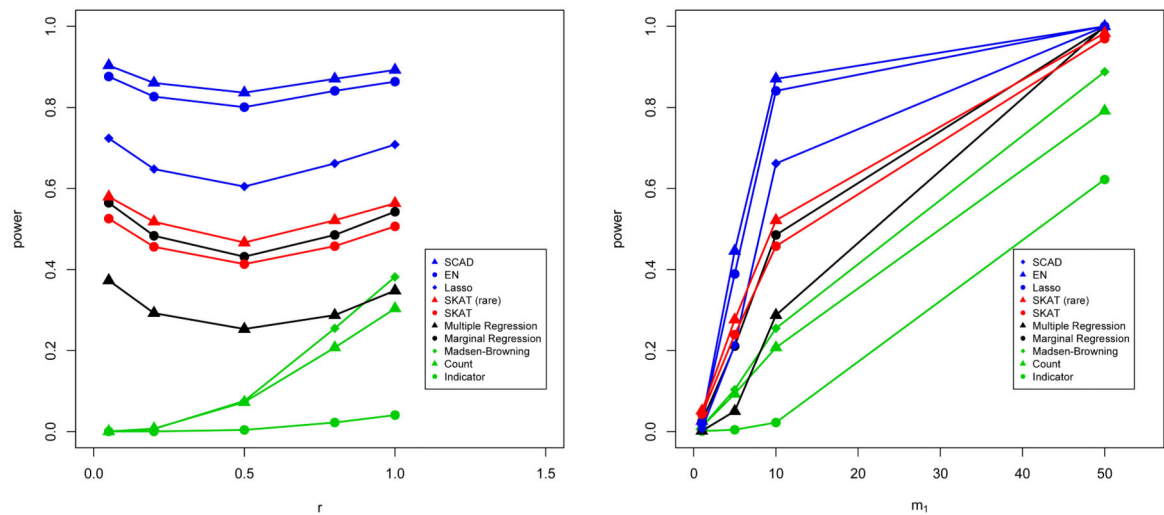
## References

- Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE, Kang HM, Marth GT, McVean GA. An integrated map of genetic variation from 1,092 human genomes. *Nature*. 2012; 491(7422):56–65. [PubMed: 23128226]
- Auer PL, Johnsen JM, Johnson AD, Logsdon BA, Lange LA, Nalls MA, Zhang G, Franceschini N, Fox K, Lange EM, Rich SS, O'Donnell CJ, Jackson RD, Wallace RB, Chen Z, Graubert TA, Wilson JG, Tang H, Lettre G, Reiner AP, Ganesh SK, Li Y. Imputation of Exome Sequence Variants into Population- Based Samples and Blood-Cell-Trait-Associated Loci in African Americans: NHLBI GO Exome Sequencing Project. *Am J Hum Genet*. 2012; 91(5):794–808. [PubMed: 23103231]
- Bacanu SA, Nelson MR, Whittaker JC. Comparison of methods and sampling designs to test for association between rare variants and quantitative traits. *Genetic Epidemiology*. 2011
- Cheung YH, Wang G, Leal SM, Wang S. A fast and noise-resilient approach to detect rare-variant associations with deep sequencing data for complex disorders. *Genet Epidemiol*. 2012; 36(7):675–85. [PubMed: 22865616]
- Cohen JC, Kiss RS, Pertsemlidis A, Marcel YL, McPherson R, Hobbs HH. Multiple rare alleles contribute to low plasma levels of HDL cholesterol. *Science*. 2004; 305(5685):869–872. [PubMed: 15297675]
- de Bakker PIW, Ferreira MAR, Jia XM, Neale BM, Raychaudhuri S, Voight BF. Practical aspects of imputation-driven meta-analysis of genome-wide association studies. *Human Molecular Genetics*. 2008; 17:R122–R128. [PubMed: 18852200]

- Dickson SP, Wang K, Krantz I, Hakonarson H, Goldstein DB. Rare variants create synthetic genome-wide associations. *PLoS Biology*. 2010; 8(1):e1000294. [PubMed: 20126254]
- Davidson, AC.; Hinkley, DV. *Bootstrap Methods and Their Applications*. Cambridge University Press; New York: 1997.
- Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, Nadeau JH. Missing heritability and strategies for finding the underlying causes of complex disease. *Nature Reviews Genetics*. 2010; 11(6):446–50.
- Epstein MP, Duncan R, Jiang Y, Conneely KN, Allen AS, Satten GA. A permutation procedure to correct for confounders in case-control studies, including tests of rare variation. *Am J Hum Genet*. 2012; 91(2):215–23. [PubMed: 22818855]
- Firmann M, Mayor V, Vidal PM, Bochud M, Pecoud A, Hayoz D, Paccaud F, Preisig M, Song KS, Yuan X, Danoff TM, Stirnadel HA, Waterworth D, Mooser V, Waeber G, Vollenweider P. The CoLaus study: a population-based study to investigate the epidemiology and genetic determinants of cardiovascular risk factors and metabolic syndrome. *Bmc Cardiovascular Disorders*. 2008; 8
- Gibson G. Hints of hidden heritability in GWAS. *Nature Genetics*. 2010; 42(7):558–60. [PubMed: 20581876]
- Gorlov IP, Gorlova OY, Sunyaev SR, Spitz MR, Amos CI. Shifting paradigm of association studies: value of rare single-nucleotide polymorphisms. *American Journal of Human Genetics*. 2008; 82(1):100–12. [PubMed: 18179889]
- Haase CL, Frikke-Schmidt R, Nordestgaard BG, Tybjaerg-Hansen A. Population-Based Resequencing of APOA1 in 10,330 Individuals: Spectrum of Genetic Variation, Phenotype, and Comparison with Extreme Phenotype Approach. *PLoS Genet*. 2012; 8(11):e1003063. [PubMed: 23209431]
- Heckman NE, Ramsay JO. Penalized regression with model-based penalties. *Canadian Journal of Statistics-Revue Canadienne De Statistique*. 2000; 28(2):241–258.
- Hesterberg T, Choi NH, Meier L, Fraley C. Least angle and  $\ell_1$  penalized regression: A review. *Statistics Surveys*. 2008; 2:61–93.
- Kyung M, Gill J, Ghosh M, Casella G. Penalized Regression, Standard Errors, and Bayesian Lasso. *Bayesian Analysis*. 2010; 5(2):369–411.
- Lee S, Wu MC, Lin X. Optimal tests for rare variant effects in sequencing association studies. *Biostatistics*. 2012
- Li BS, Leal SM. Methods for detecting associations with rare variants for common diseases: Application to analysis of sequence data. *American Journal of Human Genetics*. 2008; 83(3):311–321. [PubMed: 18691683]
- Li RQ, Li YR, Fang XD, Yang HM, Wang J, Kristiansen K. SNP detection for massively parallel whole-genome resequencing. *Genome Research*. 2009a; 19(6):1124–1132. [PubMed: 19420381]
- Li Y, Byrnes AE, Li M. To Identify Associations with Rare Variants, Just WHaIT: Weighted Haplotype and Imputation-Based Tests. *American Journal of Human Genetics*. 2010a; 87(5):728–35. [PubMed: 21055717]
- Li Y, Willer C, Sanna S, Abecasis G. Genotype imputation. *Annual Review of Genomics and Human Genetics*. 2009b; 10:387–406.
- Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genetic Epidemiology*. 2010b; 34(8):816–34. [PubMed: 21058334]
- Lin DY. An efficient Monte Carlo approach to assessing statistical significance in genomic studies. *Bioinformatics*. 2005; 21(6):781–7. [PubMed: 15454414]
- Liu EY, Buyske S, Aragaki AK, Peters U, Boerwinkle E, Carlson C, Carty C, Crawford DC, Haessler J, Hindorf LA, Marchand LL, Manolio TA, Matise T, Wang W, Kooperberg C, North KE, Li Y. Genotype Imputation of MetachipSNPs Using a Study-Specific Reference Panel of ~4,000 Haplotypes in African Americans From the Women's Health Initiative. *Genetic Epidemiology*. 2012; 36(2):107–117. [PubMed: 22851474]
- Madsen BE, Browning SR. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genetics*. 2009; 5(2):e1000384. [PubMed: 19214210]
- Maher B. Personal genomes: The case of the missing heritability. *Nature*. 2008; 456(7218):18–21. [PubMed: 18987709]

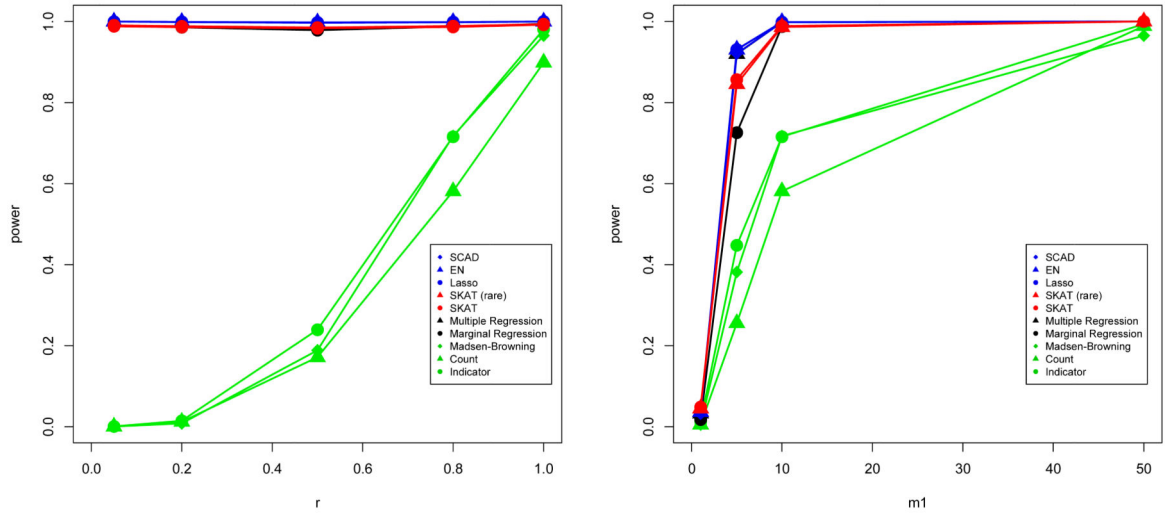
- Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, McCarthy MI, Ramos EM, Cardon LR, Chakravarti A, Cho JH, Guttmacher AE, Kong A, Kruglyak L, Mardis E, Rotimi CN, Slatkin M, Valle D, Whittemore AS, Boehnke M, Clark AG, Eichler EE, Gibson G, Haines JL, Mackay TFC, McCarroll SA, Visscher PM. Finding the missing heritability of complex diseases. *Nature*. 2009; 461(7265):747–753. [PubMed: 19812666]
- Mao X, Li Y, Liu Y, Lange L, Li M. Testing Genetic Association With Rare Variants in Admixed Populations. *Genetic Epidemiology*. 2012 n/a-n/a.
- Marchini J, Howie B. Genotype imputation for genome-wide association studies. *Nature Reviews Genetics*. 2010; 11(7):499–511.
- Morgenthaler S, Thilly WG. A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutation Research*. 2007; 615(1-2):28–56. [PubMed: 17101154]
- Neale BM, Rivas MA, Voight BF, Altshuler D, Devlin B, Orho-Melander M, Kathiresan S, Purcell SM, Roeder K, Daly MJ. Testing for an Unusual Distribution of Rare Variants. *PLoS Genetics*. 2011; 7(3):e1001322. [PubMed: 21408211]
- Nejentsev S, Walker N, Riches D, Egholm M, Todd JA. Rare variants of IFIH1, a gene implicated in antiviral responses, protect against type 1 diabetes. *Science*. 2009; 324(5925):387–9. [PubMed: 19264985]
- Nelson MR, Wegmann D, Ehm MG, Kessner D, Jean P, Verzilli C, Shen J, Tang Z, Bacanu SA, Fraser D, Warren L, Aponte J, Zawistowski M, Liu X, Zhang H, Zhang Y, Li J, Li Y, Li L, Woollard P, Topp S, Hall MD, Nangle K, Wang J, Abecasis G, Cardon LR, Zollner S, Whittaker JC, Chissole SL, Novembre J, Mooser V. An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science*. 2012; 337(6090):100–4. [PubMed: 22604722]
- Ng PC, Henikoff S. SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Research*. 2003; 31(13):3812–4. [PubMed: 12824425]
- Price AL, Kryukov GV, de Bakker PI, Purcell SM, Staples J, Wei LJ, Sunyaev SR. Pooled association tests for rare variants in exon-resequencing studies. *American Journal of Human Genetics*. 2010; 86(6):832–8. [PubMed: 20471002]
- Sanna S, Li B, Mulas A, Sidore C, Kang HM, Jackson AU, Piras MG, Usala G, Maninchedda G, Sassu A, Serra F, Palmas MA, Wood WH III, Njølstad I, Laakso M, Hveem K, Tuomilehto J, Lakka TA, Rauramaa R, Boehnke M, Cucca F, Uda M, Schlessinger D, Nagaraja R, Abecasis GR. Fine Mapping of Five Loci Associated with Low-Density Lipoprotein Cholesterol Detects Variants That Double the Explained Heritability. *PLoS Genetics*. 2011; 7(7):e1002198. [PubMed: 21829380]
- Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ, Altshuler D. Calibrating a coalescent simulation of human genome sequence variation. *Genome Research*. 2005; 15(11):1576–83. [PubMed: 16251467]
- Tibshirani R. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society (B)*. 1996; 58:267–288.
- Tzeng JY, Zhang D, Pongpanich M, Smith C, McCarthy MI, Sale MM, Worrall BB, Hsu FC, Thomas DC, Sullivan PF. Studying gene and gene-environment effects of uncommon and common variants on continuous traits: a marker-set approach using gene-trait similarity regression. *Am J Hum Genet*. 2011; 89(2):277–88. [PubMed: 21835306]
- Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. *American Journal of Human Genetics*. 2011; 89(1):82–93. [PubMed: 21737059]
- Wu TT, Lange K. Coordinate descent algorithms for lasso penalized regression. *Annals of Applied Statistics*. 2008; 2(1):224–244.
- Xie HL, Huang J. SCAD-penalized regression in high-dimensional partially linear models. *Annals of Statistics*. 2009; 37(2):673–696.
- Xu C, Ladouceur M, Dastani Z, Richards JB, Ciampi A, Greenwood CM. Multiple regression methods show great potential for rare variant association tests. *PLoS One*. 2012; 7(8):e41694. [PubMed: 22916111]

- Yi N, Liu N, Zhi D, Li J. Hierarchical Generalized Linear Models for Multiple Groups of Rare and Common Variants: Jointly Estimating Group and Individual-Variant Effects. *PLoS Genet.* 2011; 7(12):e1002382. [PubMed: 22144906]
- Zawistowski M, Gopalakrishnan S, Ding J, Li Y, Grimm S, Zollner S. Extending rare-variant testing strategies: analysis of noncoding sequence and imputed genotypes. *American Journal of Human Genetics.* 2010; 87(5):604–17. [PubMed: 21070896]
- Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B-Statistical Methodology.* 2005; 67:301–320.



### Figure 1. Power Comparison in the Absence of a Bioinformatics Tool

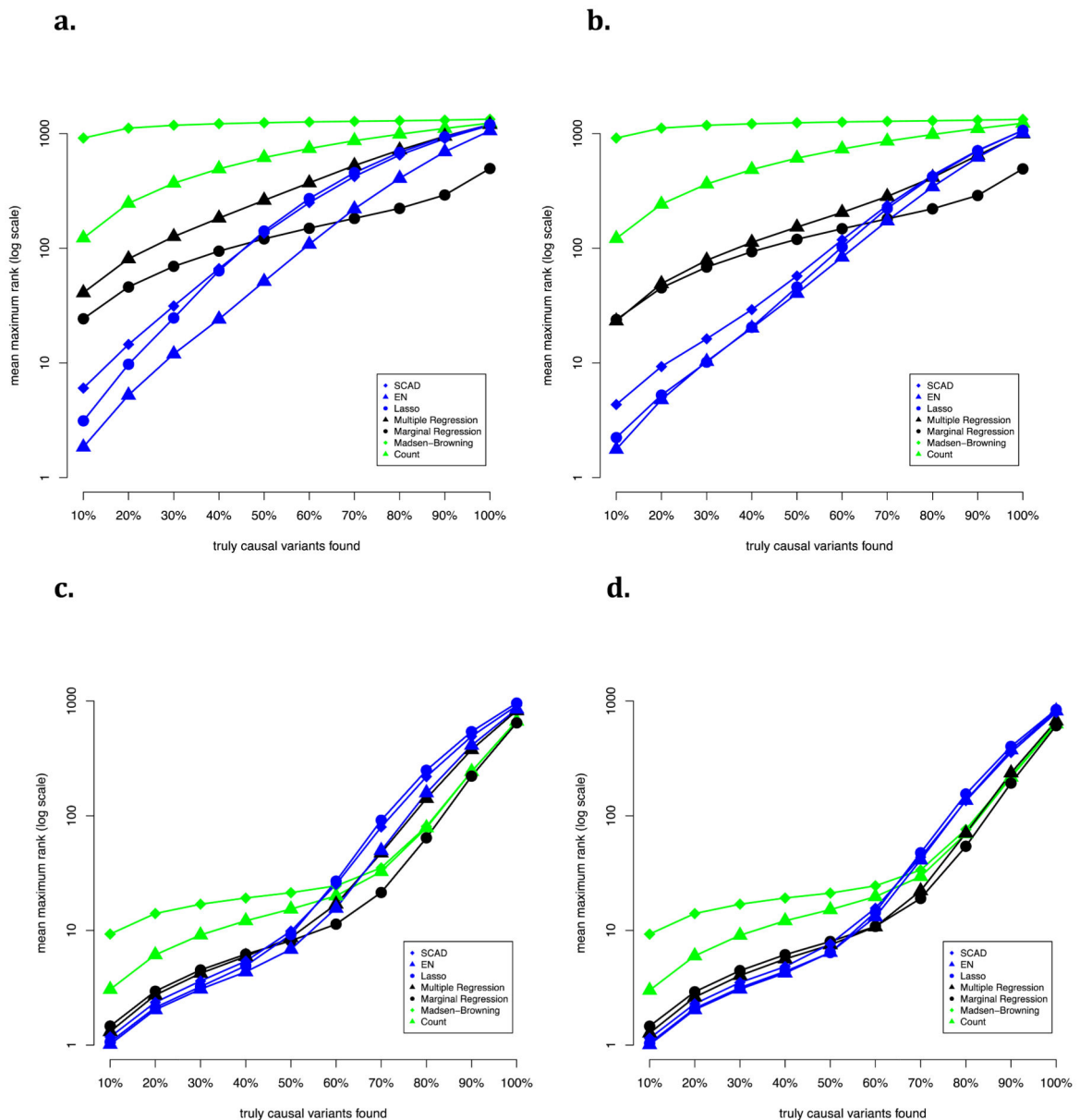
Figure 1 shows the power (Y-axis) of the different methods across a wide spectrum of  $m$  (the number of true causal variants) and  $r$  (the proportion of variants that contribute to our quantitative trait in a positive direction) in the absence of a bioinformatics tool. In Figure 1a, we fix  $m$  at 10 and show power comparisons across the entire spectrum of  $r$  (X-axis). Figure 1b shows how power changes as a function of  $m$  (X-axis) with  $r$  fixed at 0.8. Here we use the logit link function.



**Figure 2. Power Comparison in the Presence of the Good Bioinformatics Tool**

Figure 2 shows the power (Y-axis) of the different methods across a wide spectrum of  $m$  (the number of true causal variants) and  $r$  (the proportion of variants that contribute to our quantitative trait in a positive direction) in the presence of the good bioinformatics tool described in the Method section. Like in Figure 1a, we fix  $m$  at 10 and show power comparisons across the entire spectrum of  $r$  (X-axis) in Figure 2a. Similarly, Figure 2b how power of the methods changes as a function of  $m$  (X-axis) with  $r$  fixed at 0.8. Again the logit link function is used.





**Figure 3. How Far Down the Ranked List are the Truly Causal Variants when All Variants are Included?**

Figure 3a shows the number of variants that must be considered (Y-axis) in order to catch the top 10%, 20% ... 100% of truly causal variants (X-axis) in simulation when all variants are considered. We assume that the variants are ranked in order of significance. These plots aggregate true and estimated weights from all 10,000 replicates of the experiment and once again, we fix  $r$  at 0.8,  $m$  at 10 and use the logit link function. Figure 3b. takes LD buddies (variants with  $r^2 > 0.8$  with causal variant) into consideration. Figure 3c. restricts the results from 3a. to functional variants only using a good bioinformatics tool. Figure 3d. is restricted to functional variants only and takes LD buddies into account.

**Table 1**  
**Average Pearson Correlation of True and Estimated Weights ( $m=10$  and  $r=0.8$ )**

Method	All markers	Limited to functional markers
Indicator	-	-
Count	0.0126	0.2386
Madsen-Browning	0.0591	0.1225
Marginal Regression	0.1588	0.6490
Multiple Regression	0.0883	0.6537
Lasso	0.2852	0.7436
EN	0.3555	0.7787
SCAD	0.2301	0.7344
SKAT (all)	-	-
SKAT (rare only)	-	-

**Table 2**  
**Permuted P-values<sup>l</sup> on the Positive Control Gene in the Real Data Set**

Method	All variants (491)	Limited to functional variants (13)
Indicator	0.208	0.00057
Count	0.068	<b><i>0.00017</i></b> ***
Madsen-Browning	0.090	<b>0.00041</b> **
Marginal Regression	0.166	0.00420
Multiple Regression	0.136	0.00395
Lasso	<b>0.017</b> **	0.00053
EN	<b><i>0.008</i></b> ***	0.00059
SCAD	0.111	0.00078
SKAT (all)	0.329	0.00142
SKAT (rare only)	0.348	0.00142

<sup>l</sup> Except for SKAT(all) and SKAT(rare only)

\*\*\* : Most significant p-value under each column is in bold, italicized and flagged with \*\*\*.

\*\* : Second most significant p-value under each column is in bold and flagged with \*\*.