

Testing groups of genomic locations for enrichment in disease loci using linkage scan data: A method for hypothesis testing

Dimitrios Avramopoulos,^{1,2*#} Peter Zandi,^{1,3#} Adrian Gherman,² M. Daniele Fallin³ and Susan S. Bassett¹

¹ Department of Psychiatry, The Johns Hopkins University, School of Medicine, Broadway Research Building 509, 733 North Broadway, Baltimore, MD 21205, USA

² McKusick Nathans Institute for Genetic Medicine, The Johns Hopkins University, School of Medicine, 733 N. Broadway, Baltimore, MD 21205, USA

³ Department of Epidemiology, The Johns Hopkins University, Bloomberg School of Public Health, 615 N. Wolfe Street, Baltimore, MD 21205–2179, USA

#The first two authors contributed equally to the project.

*Correspondence to: Tel: +1 410 955 8323; Fax: +1 410 955 7397; E-mail: adimitr1@jhmi.edu

Date received (in revised form): 9th March 2006

Abstract

Genes for complex disorders have proven hard to find using linkage analysis. The results rarely reach the desired level of significance and researchers often have failed to replicate positive findings. There is, however, a wealth of information from other scientific approaches which enables the formation of hypotheses on groups of genes or genomic regions likely to be enriched in disease loci. Examples include genes belonging to specific pathways or producing proteins interacting with known risk factors, genes that show altered expression levels in patients or even the group of top scoring locations in a linkage study. We show here that this hypothesis of enrichment for disease loci can be tested using genome-wide linkage data, provided that these data are independent from the data used to generate the hypothesis. Our method is based on the fact that non-parametric linkage analyses are expected to show increased scores at each one of the disease loci, although this increase might not rise above the noise of stochastic variation. By using a summary statistic and calculating its empirical significance, we show that enrichment hypotheses can be tested with power higher than the power of the linkage scan data to identify individual loci. Via simulated linkage scans for a number of different models, we gain insight in the interpretation of genome scan results and test the power of our proposed method. We present an application of the method to real data from a late-onset Alzheimer's disease linkage scan as a proof of principle.

Keywords: linkage, genome scan, complex disorder, genes, group

Introduction

In complex disorders where variations in more than one gene are expected to contribute to disease risk, researchers often hypothesise that particular groups of genes or genomic locations are enriched with true disease-susceptibility genes, based on various lines of evidence. For example, groups of genes found to be differentially expressed in a case-controlled microarray experiment or the human loci syntenic to those identified by linkage in a mouse disease model are likely to be enriched in susceptibility genes. One can also hypothesise disease gene enrichment based on functional data. It can be suggested, for example, that the genes involved in glutamate

neurotransmission are enriched with schizophrenia-susceptibility genes, or — combining more than one line of evidence — that differentially expressed glutaminergic genes in particular are likely to be enriched. Researchers may wish to corroborate such hypotheses by testing whether the members of an identified group of genes are located in areas showing evidence of genetic linkage to the disease of interest. Linkage results for complex disorders are often noisy and hard to interpret, however. We propose a method for using genome-wide linkage data and the widely used non-parametric linkage (NPL) score¹ to test for the enrichment of groups of genes or genetic locations in disease-susceptibility genes.

The NPL score is designed to have, at any locus, a standard normal distribution with a mean of 0 and a standard deviation of 1 under the null hypothesis of no linkage. This means that, although for any given unlinked locus the expectation is an NPL score of 0, stochastic variation creates scores that can take positive or negative values. Under the alternative hypothesis of linkage, stochastic variation at the true disease loci is still present, but the expectation is at a value higher than 0. The magnitude of the expected value depends on the sample size, the available genetic information and the effect size of the locus. For most complex diseases, it is assumed that individual risk loci will have small effects. As a result, the superimposed stochastic variation can mask some truly linked loci or create signals where no true linkage is present, both situations leading to errors and/or failed replication studies for true disease loci. If one could have *a priori* knowledge of the true disease loci, one could achieve greater significance by studying the group of loci in concert because of the consistent trend for increased scores, even in the absence of significant scores at each one of the individual loci. As the number of true loci examined together rises, the noise from the underlying stochastic variation will asymptotically approach 0 and their average NPL score will asymptotically approach the average of their individual expectations, which is greater than 0. By contrast, for unlinked loci, the individual expectation is 0; as the number of unlinked loci examined in concert increases, their average NPL will asymptotically approach 0.

Based on these properties of the NPL score, we can use linkage analysis data to test whether a pre-defined group of loci is enriched for true disease-linked loci. This can be done by calculating the average NPL score of the group of loci and comparing it against a null distribution of average scores derived by randomly drawing groups of loci of equal size. The null hypothesis is that the proportion of true linked loci among the group of loci tested is not different from the proportion expected when choosing random loci, while the alternative hypothesis is that the proportion is greater, and hence the group is enriched. (Note that, as defined here, the proportion of true loci in the group tested corresponds to 1 minus the false discovery rate of the group.)

We assessed the power of our proposed method through simulations using a variety of disease models and varying the number of errors in location predictions. We showed that this method can be powerful, even when less than half of the examined locations are real disease loci. We must note here that we used the NPL score because it is commonly available and because its statistical properties make it easier to present our hypothesis. Since significance is determined via permutations and no distribution is assumed, however, the method is applicable to any statistic. Also, we are aware that other summary statistics, such as the product of *p*-values (more often used to show the presence of at least one true locus in a group), can serve for the same purpose. Assessment of other statistics will be the subject of future work.

Materials and methods

Generation of simulated linkage scans

We assumed a baseline risk for the disease of 0.9 per cent, representing non-genetic factors. We used the ‘-simulate’ function in the Merlin analysis package² to generate genome-wide marker data for nuclear families (two parents and four offspring), including five, ten or 20 biallelic (disease) loci carrying risk alleles of frequency *p* that independently increase the risk of disease two- or threefold (this corresponds to their relative risk). The disease allele frequency *p* was equal for all loci and was empirically adjusted to provide a population prevalence of 3 per cent. This corresponds to 70 per cent heritability (genetic/total variance), consistent with reports for many complex disorders. After generating data for a large number of families (up to 200,000 — or 1.2 million individuals), the genotypes at the disease loci for each individual were examined, the risk was determined based on those genotypes and disease status was assigned with a probability corresponding to the risk. For example, a person carrying four risk alleles with a relative risk of 2 had a probability of $2^4 \times 0.009 = 0.144$ of being affected. Sufficient families were generated every time to ascertain 1,000 sibling pairs and 60 sibling triads (1,180 sibling pairs in total) or 500 sibling pairs and 30 sibling triads (590 sibling pairs). This ratio of pairs to triads corresponds to the most efficient choice, given the observed simulated families, but it is not far from common sibling size distributions in the complex disease literature. The genotypes of all markers, excluding the biallelic disease loci, were used for genome scans for linkage using the Merlin software² to calculate NPL scores across the genome. Markers other than the disease loci had six equiprobable alleles, spaced 10 centimorgans (cM) apart, and there were no missing data. A total of 359 microsatellite markers were simulated, starting on each chromosome at genetic position 0 and placing one marker every 10 cM until the end of the chromosome; therefore, a genetic length of 0 to less than 10 cM was at the telomere of each chromosome. All families had a size of six, with two parents and four offspring. Although this might be a slightly large family size compared with today's average in the Western world, it is less so for families that are ascertained today through their adult affected offspring. This family size was also necessary to make the generation of enough pedigrees for ascertainment computationally feasible. For each simulated scan, the location of the disease loci varied and these were placed randomly in the genome, allowing for co-localisation of more than one disease locus and for zero distance with scan markers if it so happened by chance.

Models examined

In order to assess the power of our method under many possible scenarios, we studied multiple genome scans under multiple disease parameters including: 1) The number of disease

genes was set to five, ten or 20; 2) The increase in risk from each risk allele was set to 2 or 3 relative to the baseline risk; 3) The number of ascertained families was 500 pairs + 30 triplets or 1,000 pairs + 60 triplets. For each of the 12 possible sets of parameters, 25 genome scans were generated. For each genome scan, we examined a number of different scenarios regarding the number of true and non-true disease loci in the group to be tested. When some non-true locations, or not all the real locations, were included, the groups were chosen 100 times at random to account for the stochastic variation inherent to the selection. In order to determine significance for each of the 100 selected groups, the average NPL of each was compared with the null distribution formed by the average scores of 1,000 random groups of equal size — that is, groups chosen without taking into account whether or not the included loci correspond to disease gene locations. The 2,500 empirical significance values obtained from 25 scans \times 100 group permutations were used to determine the power for each model (each cell on Table 1). Since in these simulated data the alternative hypothesis (as stated above) is always true, the number of times that the empirical significance is less than the desired significance level, α , corresponds to the power.

Results

Table 2 summarises our general observations from the simulated genome scans with 1,180 affected sibling pairs. A linkage peak was considered to contain a real disease locus if the NPL scores between the original location of the binary disease marker and the observed peak did not drop by more than one unit less than the score at the peak. Even with five loci of relative risk 3, the top NPL score did not reach genome-wide significance on most scans (25 scan average = 4.22), according to the criteria proposed by Lander and Kruglyak³ (for the genome-wide significant $p = 2.2 \times 10^{-5}$, an NPL score of 4.4 is required), in accordance to what has been observed in real data analyses⁴ and predicted by Risch and Merikangas.⁵ It is notable and encouraging, however, that, across the models we tested, 40–92 per cent of scans showed strongest linkage at a real locus. We also observed that even when there are only five true loci on average, one of these loci is not among the top ten peaks of a scan and would therefore not be detected. Again, this observation is very much in agreement with the experience from linkage studies of complex disorders, as many linkage findings that have been considered to carry strong evidence have often not been replicated in subsequent studies of different pedigrees. As expected, the NPL scores and the fraction of true positives among the top linkage peaks decrease as the number of disease loci increases and as their relative risk decreases. An increase in the fraction of true findings is counter-intuitively observed as the number of true disease loci and the number of top linkage peaks examined increases to 20;

however, this does not correspond to an increase in the fraction of real genes identified. When looking at Table 2, the reader should keep in mind that when there are only five real genes and a set of 20 loci is tested, the maximum possible fraction of true loci in the set is $5/20 = 25$ per cent. Overall, our observations confirm that our confidence in the linkage peaks of a single scan should be somewhat reserved until we observe replication, but also that non-replication of a linkage finding does not necessarily discredit a positive finding. In other words, it will take more than a few linkage scans to develop strong confidence in the location of true susceptibility loci for a complex disorder.

Table 1 presents an evaluation of the power of the approach we propose here for examining multiple genomic locations for linkage using a summary statistic, namely the average NPL score. In particular, it shows the power to detect enrichment by examining the average NPL score at levels of $\alpha = 0.05$ and 0.01 and for different disease models calculated through our computer simulations. For a 1,180 sibling pair scan, and at the nominal level of $\alpha = 0.05$, it is of interest that for a relative risk of 3 and for as many as ten disease loci, we can observe significance with power > 80 per cent even if only one-third of the locations in the group are true. For a relative risk of 2, we have 80 per cent power if half of the locations in the group are true. For 20 segregating loci and a relative risk of 3, we can only tolerate ten non-real locations in a group that includes all 20 correct locations if we wish to have 80 per cent power. As expected, the power is reduced with smaller sample sizes. Figure 1 provides a three-dimensional graph showing how power increases when there are fewer real loci contributing to the risk and when more of these are included in the group. As the group gets larger and the fraction of true loci that are included is reduced, however, the power decreases.

Based on our observations on the true positive enrichment of the top linkage peaks (Table 2) and on the power of averaging NPL scores (Table 1), we decided to test our approach on real data. Our simulations indicate that, as expected, the group of top findings of a linkage scan is enriched in true disease locations. Therefore, provided that there are not too many true loci with too small effects, when this group of top linkage peaks is tested against an independent linkage scan, it should show a significantly elevated average NPL score. We used our genome scan data on late-onset Alzheimer's disease (LOAD) to test this hypothesis. This scan was performed on a previously described collection of pedigrees from the National Institute of Mental Health genetics initiative,⁶ for which we have previously reported genome scan results.⁷ The one gene known to be involved in LOAD is *APOE*,⁸ and its behaviour in terms of risk is very similar to our simulated models.^{9,10} It has been suggested that another 4–5 loci with effects similar to *APOE* may be involved in LOAD.¹¹ Adopting a study design that enabled us to keep most variables equal and yet have two independent scans, we sorted the pedigrees by their assigned

Table 1. The power of our method for different simulated models (five, ten or 20 disease loci, 1,180 or 590 sibling pairs, relative risk (RR) of 2 or 3) different levels of enrichment for true loci and different levels of significance.

| # Real loci | Group composition | | Power: 1,180 sibling pairs, RR = 3% | | Power: 1,180 sibling pairs, RR = 2% | | Power: 590 sibling pairs, RR = 3% | | Power: 590 sibling pairs, RR = 2% | |
|-------------|-------------------|-------|-------------------------------------|-----------------|-------------------------------------|-----------------|-----------------------------------|-----------------|-----------------------------------|-----------------|
| | Real loci | False | $\alpha = 0.05$ | $\alpha = 0.01$ | $\alpha = 0.05$ | $\alpha = 0.01$ | $\alpha = 0.05$ | $\alpha = 0.01$ | $\alpha = 0.05$ | $\alpha = 0.01$ |
| | 5 | 0 | 100 | 100 | 100 | 92 | 100 | 96 | 100 | 84 |
| | 5 | 3 | 100 | 98 | 90 | 71 | 94 | 88 | 80 | 51 |
| | 5 | 5 | 99 | 94 | 81 | 58 | 93 | 79 | 70 | 40 |
| 5 | 5 | 10 | 95 | 76 | 69 | 42 | 84 | 58 | 53 | 25 |
| | 3 | 0 | 100 | 95 | 88 | 72 | 97 | 88 | 82 | 52 |
| | 3 | 2 | 94 | 78 | 71 | 43 | 84 | 61 | 57 | 30 |
| | 3 | 5 | 80 | 52 | 54 | 27 | 67 | 37 | 38 | 16 |
| | 10 | 0 | 100 | 100 | 100 | 92 | 100 | 92 | 72 | 48 |
| | 10 | 5 | 99 | 92 | 89 | 63 | 88 | 64 | 54 | 29 |
| | 10 | 10 | 94 | 77 | 78 | 48 | 75 | 48 | 46 | 22 |
| | 10 | 20 | 80 | 54 | 58 | 30 | 60 | 30 | 35 | 14 |
| 10 | 5 | 0 | 90 | 72 | 75 | 48 | 77 | 44 | 48 | 20 |
| | 5 | 3 | 75 | 49 | 56 | 28 | 53 | 26 | 32 | 13 |
| | 5 | 5 | 68 | 42 | 49 | 24 | 47 | 21 | 29 | 11 |
| | 5 | 10 | 53 | 28 | 35 | 15 | 35 | 16 | 22 | 7 |
| | 20 | 0 | 100 | 76 | 68 | 44 | 65 | 35 | 32 | 12 |
| | 20 | 10 | 80 | 49 | 56 | 31 | 47 | 23 | 26 | 10 |
| | 20 | 20 | 67 | 37 | 47 | 21 | 40 | 18 | 24 | 10 |
| | 20 | 40 | 49 | 24 | 36 | 16 | 31 | 14 | 17 | 6 |
| | 15 | 0 | 84 | 58 | 61 | 35 | 53 | 26 | 29 | 13 |
| 20 | 15 | 10 | 64 | 34 | 45 | 20 | 40 | 16 | 23 | 9 |
| | 15 | 20 | 51 | 23 | 36 | 15 | 29 | 13 | 18 | 7 |
| | 15 | 30 | 41 | 19 | 29 | 12 | 26 | 9 | 14 | 5 |
| | 10 | 0 | 67 | 40 | 47 | 23 | 41 | 18 | 25 | 9 |
| | 10 | 5 | 51 | 24 | 38 | 15 | 32 | 13 | 19 | 6 |
| | 10 | 10 | 41 | 19 | 31 | 11 | 26 | 9 | 16 | 6 |
| | 10 | 15 | 35 | 14 | 24 | 9 | 23 | 8 | 15 | 5 |

identification numbers (signifying collection site and collection sequence) and split them at the point that gives two sample sets (sets A and B) of equal numbers of sibling pairs (296 pairs each). We then ran a genome-wide linkage

analysis on both sets, ranked the top scoring 30 locations from scan A, selected groups of five, ten, 15...30 locations starting from the top and tested their average NPL in scan B. We note that splitting the data has no benefit for gene discovery, but

Table 2. Results of simulated scans with 1,180 sibling pairs regarding their success in identifying the disease gene locations.

| Simulated model | # Disease loci | Number of top findings | | | | |
|---------------------------|----------------|-----------------------------------|-----------------------------------|-----------------------------------|------------------------------------|------------------------------------|
| | | 1 Real loci (%) Average NPL | 3 Real loci (%) Average NPL | 5 Real loci (%) Average NPL | 10 Real loci (%) Average NPL | 20 Real loci (%) Average NPL |
| H = 0.7, K ≈ 0.03, RR = 3 | 5 | 0.92 (92%) | 2.28 (76%) | 3.16 (63%) | 4.04 (40%) | 4.4 (22%) |
| | | 4.22 | 3.70 | 3.36 | 2.86 | 2.29 |
| | 10 | 0.72 (72%) | 2.12 (71%) | 3.24 (65%) | 4.8 (48%) | 6.32 (32%) |
| | | 3.55 | 3.22 | 2.98 | 2.58 | 2.12 |
| | 20 | 0.52 (52%) | 1.48 (49%) | 2.24 (45%) | 4.24 (42%) | 8.04 (40%) |
| | | 3.15 | 2.87 | 2.71 | 2.41 | 2.03 |
| H = 0.7, K ≈ 0.03, RR = 2 | 5 | 0.75 (75%) | 1.54 (51%) | 1.96 (39%) | 2.79 (28%) | 3.58 (18%) |
| | | 3.24 | 2.91 | 2.68 | 2.34 | 1.93 |
| | 10 | 0.56 (56%) | 1.44 (48%) | 2.16 (43%) | 3.32 (33%) | 5.08 (25%) |
| | | 3.10 | 2.83 | 2.62 | 2.30 | 1.91 |
| | 20 | 0.4 (40%) | 1 (33%) | 1.64 (33%) | 3.24 (32%) | 5.56 (28%) |
| | | 2.79 | 2.54 | 2.39 | 2.12 | 1.76 |

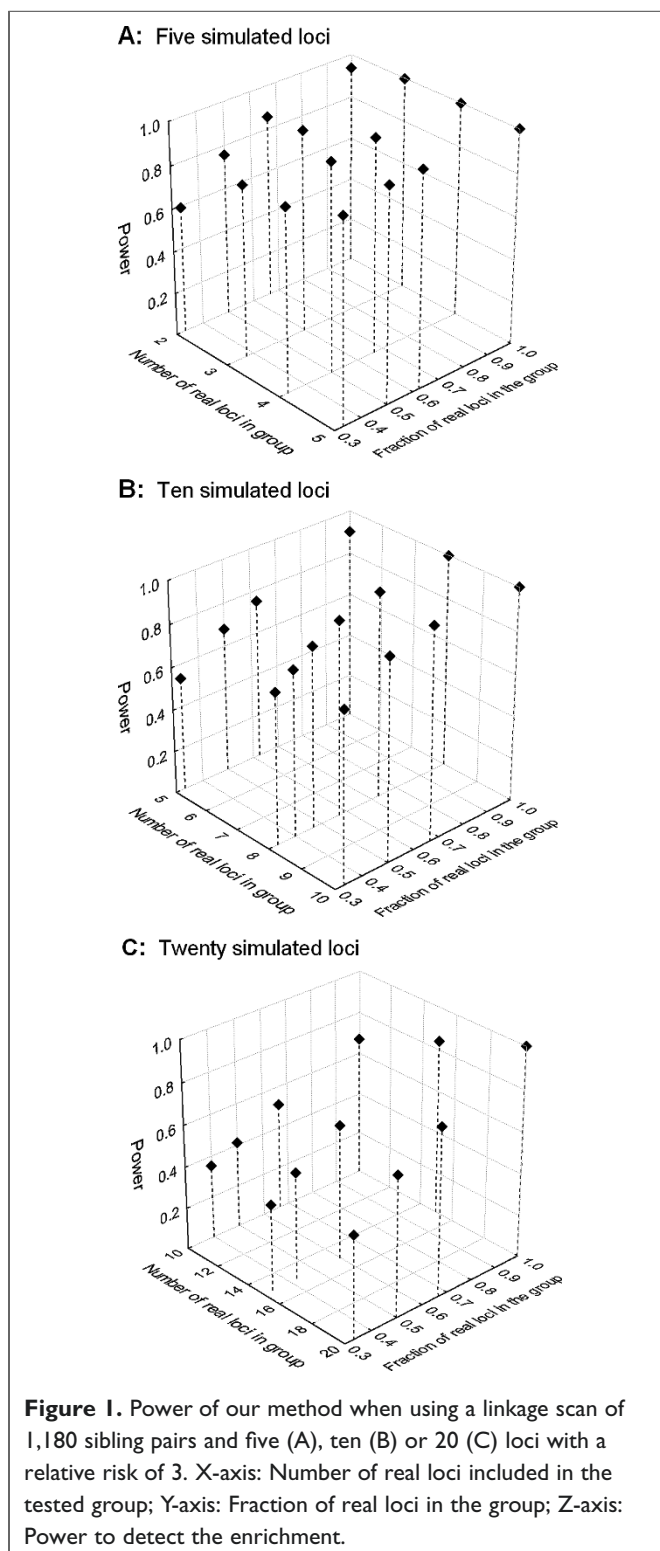
Simulation parameters: H = heritability, K = prevalence, RR = relative risk for each risk allele. The number of simulated disease loci is shown. For the one, three, five, ten and 20 top scoring loci for each of the 25 simulated genome scans, we show how many coincided with true disease loci (and their percentages), as well as their average non-parametric linkage (NPL) score.

we did this here as an exercise to show proof of principle because it provided us with a hypothesis that we could readily test using our method, namely that the peaks of a genome scan are enriched in true loci. Given the small sample sizes (296 sibling pairs per scan), our power might have been low, since the underlying model is unknown; however, we viewed this analysis as exploratory. Table 3 shows the empirical significance obtained by selecting the best five and up to 30 locations from the top linkage peaks based on scan A, and testing their average NPL against the data from scan B. Although the mean NPL of the five top locations was not significantly high, once the number was raised to ten and 15, the scores were significant, suggesting enrichment in disease gene locations. We consider that this not only validates our method but that it is also very encouraging regarding the validity of the findings of our genome scan, suggesting that the top linkage peaks are indeed enriched in real disease loci to a significant degree. Although some might consider this notion to be obvious, it is contingent on the underlying disease model and might not necessarily be true. Based on the observations from the simulated linkage scans, and the expected low power of this test on 296 sibling pairs, this also suggests that the number of substantial disease loci is not too great and that their relative risks are not too small. As we performed comparisons in six groups, we next wanted to

see if our findings were significant at the study-wide level. The strong correlation between the tested groups makes Bonferroni correction too conservative, so we tested this empirically. We chose 10,000 groups of 30 loci and tested inclusive subgroups of five, ten, 15...30 members, as we did with the real data against the scan B results. A p -value of 0.016 or smaller in any sub-group was obtained 606 times, providing a study-wide significance of 0.06.

Discussion

We have shown how one can test for the enrichment of a group of genomic locations for disease loci using linkage genome scan data. Candidate groups of genomic locations can arise from multiple types of data. For example, one can compare the results of two independent genome scans, as described here for Alzheimer's disease, in the same or different organisms. Alternatively, one can test prior results of expression studies or genome-wide association analyses, or genes belonging to specific pathways or interacting with a suspected disease gene. The method could be extended to applying weights to individual locations based on the strength of prior evidence. This is highly intuitive for testing locations that carry a score or a significance value (such as



linkage, expression or association results) but less so for other types of groups (interacting proteins, members of a functional group, etc). One could also extend the approach by testing groups on data other than linkage results, yet such approaches

Table 3. Application of our method to real data. Two sets of pedigrees were used for scans A and B. Groups of top linkage peaks from scan A (their size is shown in column 1) were then tested for enrichment on the results of scan B. Column 2 shows the empirical p -values for these groups.

| Top locations from scan A | Significance on scan B |
|---------------------------|------------------------|
| 5 | 0.220 |
| 10 | 0.016 |
| 15 | 0.048 |
| 20 | 0.077 |
| 25 | 0.089 |
| 30 | 0.094 |

require further method development because there can be a number of issues that need to be addressed.

The use of sum statistics for SNP association data has been described previously by Wille *et al.*,¹² Hoh *et al.*¹³ and Kim *et al.*¹⁴ The goals of these investigators, however, were different to ours. These authors sought methods to test for associations in multilocus disorders, with the notion of increasing power to detect associations with any one of the loci by examining groups of SNPs or other DNA markers in concert. By contrast, we sought to develop a method specifically for testing the hypothesis that a group of pre-selected genomic locations is enriched for disease loci. Our method is suitable for testing any group of genes or locations on pre-existing data. In fact, our method could complement and add to the validity of the findings from other SNP set association studies.

There is one important pitfall about which investigators need to be cautious. It is necessary to make sure that the linkage data used for testing the enrichment hypothesis were not in any way used to generate the hypothesis. For example, if one tests genes that have been reported to be associated with a disease, it is necessary to use linkage data generated and/or published after the associations, as there is a strong bias towards association testing in linked regions. If the linkage data were known before the association studies, the genes might have been examined because of the positive linkage scores and testing their scores on the same linkage scan is certain to give a false positive result. For example, there are numerous association studies on Alzheimer's disease and we could have used our linkage data to test whether the group of positive findings is enriched for true genes. The pedigrees used in our study, however, have been publicly available and used for genome scans since 1999.¹⁵ Many association studies that followed were biased towards examining linked regions and thus the positive findings would have a similar bias. We would need results from an unbiased genome screen for association to perform a valid test for enrichment. One

also needs to consider that although the power of the method is substantial, it will quickly diminish if multiple hypotheses of little merit are examined, as this will require substantial correction for multiple comparisons. Additionally, as the true underlying disease model is not known, negative results cannot be taken as evidence against a hypothesis and must be interpreted with caution. Although failing to reject the null hypothesis might suggest that the alternative is wrong, it might also be due to decreased power resulting from the small effect of individual genes, the large number of genes involved, insufficient enrichment of the selected locations in true disease genes or the small number of pedigrees in the linkage study. Regarding the last point, the approach could be extended to simultaneous examination of two or more linkage scans to increase power without the need to combine the genotype data with all the inherent difficulties of doing so. One can simply perform permutations of the same group of random loci on both scans and examine the distribution of the combined average NPL score against the observed average of the two scans for the tested group.

Our simulation data can provide some guidance on the optimal selection of group size. As Figure 1 shows, when less than 50 per cent of the loci in the group are real, the power starts to diminish. Significant loss of power is also observed when less than half of all true loci are included in the group (Table 1); thus, we suggest using the maximum group size that does not exceed twice the predicted number of true loci. Our data on LOAD support this, as the predicted number of loci conferring a relative risk of 2–3 is five,¹¹ and we obtained our strongest finding with a group of size of ten. When information on an expected number of disease genes is available, we suggest avoiding multiple comparisons by defining *a priori* the size of the tested group to roughly twice that number. If one wishes to test multiple group sizes, correction for the multiple correlated comparisons is required using empirical methods. As we observed in our example in Alzheimer's disease, the predicted group of ten loci would have provided the highest significance, while testing six groups resulted in a study-wide *p*-value of 0.06. Variations in group size can be useful in determining the most enriched group, yet it might be best to perform this analysis after significance has been established.

Our example using Alzheimer's disease linkage data showed how positive findings can not only confirm a hypothesis — in this case, confirm that a significant proportion of disease loci are amidst the top linkage findings — but also lead to insight regarding the possible underlying model. It has been previously proposed that about five loci, each conferring a relative risk of 2–3 for LOAD, segregate in the population.¹¹ According to Table 2, for 1,180 sibling pairs and five loci with a relative risk of 3, we would expect that 3.2 of the top five and four of the top ten linkage peaks would be real. These numbers would be 2.5 and 3.4, respectively, for 590 sibling pairs. If we compared these against a linkage scan of 590

sibling pairs, extrapolating from Table 1, we would expect to have somewhat more than 80 per cent power to detect this degree of enrichment. Although our sample for both genome scans was about half the size of this sample and presumably had significantly less power, we detected the enrichment in our data. Having a positive finding in this analysis that is consistent with the proposed number of loci and relative risks is very encouraging, as it suggests that linkage analysis has pointed to some truly linked regions in our Alzheimer's genome scan.

The analytical approach we propose here is simple and, since it calculates the significance of findings based on permutations, robust to type I errors, provided that the prediction of the genomic locations to be grouped and analysed is in no way biased by the linkage data on which the test will be performed. We showed that the approach has substantial power under disease models with a moderate number of risk genes and moderate relative risks. We believe that as more and more diverse data accumulate through the various high-throughput technologies, it is increasingly important to devise more methods of combining and cross-validating the resulting information that will help us succeed in our effort to understand complex disorders.

Acknowledgments

This work was supported by grants from the N.I.A. to D.A. (AG022099) and to S.S.B. (AG021804), and by an award from the Neurosciences Education and Research Foundation to D.A.

References

1. Kruglyak, L., Daly, M.J., Reeve-Daly, M.P. and Lander, E.S. (1996), 'Parametric and nonparametric linkage analysis: A unified multipoint approach', *Am. J. Hum. Genet.* Vol. 58, pp. 1347–1363.
2. Abecasis, G.R., Cherny, S.S., Cookson, W.O. and Cardon, L.R. (2002), 'Merlin — Rapid analysis of dense genetic maps using sparse gene flow trees', *Nat. Genet.* Vol. 30, pp. 97–101.
3. Lander, E. and Kruglyak, L. (1995), 'Genetic dissection of complex traits: Guidelines for interpreting and reporting linkage results', *Nat. Genet.* Vol. 11, pp. 241–247.
4. Altmuller, J., Palmer, L.J., Fischer, G. *et al.* (2001), 'Genomewide scans of complex human diseases: True linkage is hard to find', *Am. J. Hum. Genet.* Vol. 69, pp. 936–950.
5. Risch, N. and Merikangas, K. (1996), 'The future of genetic studies of complex human diseases', *Science* Vol. 273, pp. 1516–1517.
6. Blacker, D., Haines, J.L., Rodes, L. *et al.* (1997), 'ApoE-4 and age at onset of Alzheimer's disease: The NIMH genetics initiative', *Neurology* Vol. 48, pp. 139–147.
7. Blacker, D., Bertram, L., Saunders, A.J. *et al.* (2003), 'Results of a high-resolution genome screen of 437 Alzheimer's disease families', *Hum. Mol. Genet.* Vol. 12, pp. 23–32.
8. Strittmatter, W.J., Saunders, A.M., Schmechel, D. *et al.* (1993), 'Apolipoprotein E: High-avidity binding to beta-amyloid and increased frequency of type 4 allele in late-onset familial Alzheimer disease', *Proc. Natl. Acad. Sci. USA* Vol. 90, pp. 1977–1981.
9. Evans, D.A., Bennett, D.A., Wilson, R.S. *et al.* (2003), 'Incidence of Alzheimer disease in a biracial urban community: Relation to apolipoprotein E allele status', *Arch. Neurol.* Vol. 60, pp. 185–189.

10. Tang, M.X., Stern, Y., Marder, K. *et al.* (1998), 'The APOE-epsilon4 allele and the risk of Alzheimer disease among African Americans, whites, and Hispanics', *JAMA* Vol. 279, pp. 751–755.
11. Warwick Daw, E., Payami, H., Nemens, E.J. *et al.* (2000), 'The number of trait loci in late-onset Alzheimer disease', *Am. J. Hum. Genet.* Vol. 66, pp. 196–204.
12. Wille, A., Hoh, J. and Ott, J. (2003), 'Sum statistics for the joint detection of multiple disease loci in case-control association studies with SNP markers', *Genet. Epidemiol.* Vol. 25, pp. 350–359.
13. Hoh, J., Wille, A. and Ott, J. (2001), 'Trimming, weighting, and grouping SNPs in human case-control association studies', *Genome Res.* Vol. 11, pp. 2115–2119.
14. Kim, S., Zhang, K. and Sun, F. (2003), 'Detecting susceptibility genes in case-control studies using set association', *BMC Genet.* Vol. 4, p. S9.
15. Kehoe, P., Wavrant-De Vrieze, F., Crook, R. *et al.* (1999), 'A full genome scan for late onset Alzheimer's disease', *Hum. Mol. Genet.* Vol. 8, pp. 237–245.