

Power comparison of different methods to detect genetic effects and gene-environment interactions

Rémi Kazma*^{1,2}, Marie-Hélène Dizier^{2,1}, Michel Guilloud-Bataille^{2,1}, Catherine Bonaïti-Pellié^{2,1} and Emmanuelle Génin^{2,1}

Address: ¹Université Paris-Sud, UMR-S 535, Villejuif, 94817, France and ²INSERM UMR-S 535, Villejuif, 94817, France

Email: Rémi Kazma* - kazma@vjf.inserm.fr; Marie-Hélène Dizier - dizier@vjf.inserm.fr; Michel Guilloud-Bataille - guilloud@vjf.inserm.fr; Catherine Bonaïti-Pellié - bonaiti@vjf.inserm.fr; Emmanuelle Génin - genin@vjf.inserm.fr

* Corresponding author

from Genetic Analysis Workshop 15
St. Pete Beach, Florida, USA. 11–15 November 2006

Published: 18 December 2007

BMC Proceedings 2007, 1(Suppl 1):S74

This article is available from: <http://www.biomedcentral.com/1753-6561/1/S1/S74>

© 2007 Kazma et al; licensee BioMed Central Ltd.

This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Identifying gene-environment ($G \times E$) interactions has become a crucial issue in the past decades. Different methods have been proposed to test for $G \times E$ interactions in the framework of linkage or association testing. However, their respective performances have rarely been compared. Using Genetic Analysis Workshop 15 simulated data, we compared the power of four methods: one based on affected sib pairs that tests for linkage and interaction (the mean interaction test) and three methods that test for association and/or interaction: a case-control test, a case-only test, and a log-linear approach based on case-parent trios. Results show that for the particular model of interaction between tobacco use and Locus B simulated here, the mean interaction test has poor power to detect either the genetic effect or the interaction. The association studies, i.e., the log-linear-modeling approach and the case-control method, are more powerful to detect the genetic effect (power of 78% and 95%, respectively) and taking into account interaction moderately increases the power (increase of 9% and 3%, respectively). The case-only design exhibits a 95% power to detect $G \times E$ interaction but the type I error rate is increased.

Background

Gene-environment ($G \times E$) interactions are likely to play an important role in multifactorial diseases. The detection of $G \times E$ interaction can be of major interest in epidemiological studies to help identify subgroups of the population that are at high risk of disease and at which prevention and screening programs should be targeted. The presence of interaction can conceal environmental

and/or genetic effects if not considered in the analysis [1]. On the other hand, taking it into account may either enhance or reduce the power to detect genetic susceptibility factors, depending on the parameters inherent to the model underlying disease susceptibility [2,3]. With this in mind, many statistical methods have been developed in the past decades, either to directly investigate $G \times E$ interaction or to enhance detection of genetic factors by taking

into account exposure status. They can be classified according to the design followed, the kind of data used, and the hypothesis tested [1,4].

The purpose of our work is to compare the power of different methods to detect the effect of Locus B and its interaction with history of tobacco use. We used the simulated data (Problem 3) of Genetic Analysis Workshop 15 (GAW15) with knowledge of the "answers" and compared four methods to test for genetic effect and/or $G \times E$ interaction. The first method, referred to as the mean interaction test (MIT) method [5], tests linkage and $G \times E$ interaction among sib pairs. It is compared to three association testing methods: a log-linear-modeling approach [6] that uses case-parent triads and a case-control design [4], both of which test for the effect of the gene and $G \times E$ interaction, and a case-only design [7] that tests for interaction only.

Methods

One hundred replicates were studied at the disease susceptibility Locus B that controls the effect of smoking on rheumatoid arthritis risk. In each replicate, 1500 affected sib pairs were considered for the MIT, 1500 case-parent trios for the log-linear method, 1500 cases and controls for the case-control design, and only the 1500 cases for the case-only test. We also studied smaller sample sizes (500 trios and 750 cases and controls) in order to compare the three association methods for the same number of genotyped individuals. Cases were obtained by considering the first affected case in each sib-pair and controls were the first 1500 control subjects among the 2000 available for each replicate.

Because none of the single-nucleotide polymorphisms (SNPs) close to Locus B were in linkage disequilibrium with this locus, we used genotypes of all the individuals at that locus for association tests and the exact identity-by-descent (IBD) provided in the Problem 3 "answers" for the linkage test. For the exposure status, we considered the lifetime smoking status and did not account for the indirectly increased risk through smoke effect on IgM.

The four following methods were compared.

Mean interaction test

The MIT developed by Gauderman and Siegmund [5] is an extension of the mean sharing test [8] to account for $G \times E$ interaction. It compares the proportion of alleles shared IBD, π , which is expected to be equal to 0.5 under the null hypothesis of no linkage, across the three groups of affected sib pairs differing for the number of exposed sibs (2, 1, or 0). The following regression model is used: $\pi_i = \pi + \beta(X_i - \bar{X}) + \varepsilon_i$, where π is the intercept and β the regression coefficient for the exposure, with X_i the covari-

ate of exposure centered on its mean \bar{X} . We conducted analysis using the coding scheme consisting of two variables (X_{EE} and X_{EU}) contrasting sib pairs with 2, 1, or 0 exposed sibs. The null hypothesis of no linkage is tested by the likelihood ratio test (LRT): $T_{\pi\beta} = 2[\ln\{L(\pi = 0.5, \beta = 0)\} - \ln\{L(\pi, \beta)\}]$, which follows a 50:50 mixture distribution of two and three degrees of freedom (df) χ^2 . The alternative hypothesis corresponds to linkage with or without $G \times E$ interaction.

In its original presentation, the mean interaction test method allows accounting for $G \times E$ interaction in the search for linkage but does not test for $G \times E$ interaction. We therefore developed a LRT for $G \times E$ interaction: $T_{\beta} = 2[\ln\{L(\pi, \beta = 0)\} - \ln\{L(\pi, \beta)\}]$. This test follows a 2 df χ^2 distribution.

Log-linear-modeling approach for case-parent triads

Proposed by Umbach and Weinberg [6], this method consists of comparing the conditional genotype distribution of exposed cases, given parental genotypes, versus that of unexposed cases. Briefly, case-parent triads are divided into 20 categories based on the parental genotypes, the genotype of the case, and the exposure status of the case. The expected number of triads can be expressed according to a log-linear model [3,6]. LRT are performed to test for 1) a gene effect ignoring $G \times E$ interaction (which follows a 2 df χ^2), 2) a gene effect accounting for $G \times E$ interaction (which follows a 4 df χ^2), and 3) a $G \times E$ interaction (which follows a 2 df χ^2). Fit of the data with a dominant model is also tested as the true model was dominant.

Case-control design

Case-control designs have been widely used to compare risks of developing a disease according to their genotype and exposure status [4]. Odds-ratios (OR) associated with the exposure, the genotypes, and their interaction factors are estimated and tested for significance. Three likelihood ratio tests are performed: a 2 df χ^2 test for genetic effect alone, a 4 df χ^2 test for genetic effect accounting for $G \times E$ interaction, and 2 df χ^2 test of $G \times E$ interaction. Fit of the data with a dominant model is tested using a 2 df LRT.

Case-only design

Case-only studies [4,7] test the interaction between an exposure and a genotype among case subjects only. This type of design assesses the departure from a multiplicative scale, assuming independence between both factors. To test for the interaction, a 2 df LRT of homogeneity between the genotype distribution in exposed and unexposed cases is performed.

Powers of the different tests were estimated by determining the number of replicates among the 100 replicates that were significant at a nominal 0.05 type I error rate. Type I

error rates to test for $G \times E$ interaction are estimated on the seven loci (A, C-H) that are not supposed to interact with lifetime smoking status.

Results

Table 1 gives the mean proportion of alleles shared IBD in the whole sample, and in each of the three sib-pair categories of exposure. Table 2 shows the power of the different tests. We found that MIT has almost no power to detect linkage even when accounting for $G \times E$ interaction. This could have been expected given the proportion of alleles shared IBD in the whole sample and in each of the three categories based on exposure. Indeed, these proportions are very close to the null expectation of 0.5 (Table 1).

With the log-linear model, the power to detect the gene effect is 78% and is increased to 87% when accounting for $G \times E$ interaction. Thus, there is a gain in power to detect the gene effect when accounting for $G \times E$ interaction under the simulated model. For the case-control design, the power to detect the gene effect is 95% and improves to 98% when accounting for interaction. As shown in Figure 1, the p -values of test accounting for $G \times E$ are smaller than those of the test not accounting for $G \times E$ for most of the replicates and similar trends (gain or loss of power) are observed between the two methods in 74% of the replicates.

Concerning the detection of the $G \times E$ interaction, we found that the case-only design is by far the most powerful test. It reaches 95% power; the case-control design only reaches 69%, the log-linear approach, 53%; and the linkage test (MIT), 12%. When constraining the number of genotyped individuals to be the same for the three association methods, the differences in power are even more pronounced. Figure 2 shows the p -values of the $G \times E$ interaction test for the log-linear-modeling approach, the case-control, and the case-only designs for the first 25 replicates. We observe that it is generally in the same replicates that the different methods give the most significant results, with the highest significance achieved for the case-only method.

Estimates of interaction factors presented in Table 2 do not seem to comply with a dominant model, and indeed a dominant model is rejected in 60% of the replicates with the case-control and in 46% of the replicates with the log-linear model.

Average type I error rates for the interaction test over the seven loci were 13% for the case-only design (ranging from 5% for Locus H to 26% for Locus C), 10% for the case-control (from 4% for Locus H to 30% for Locus C), and 8% (ranging from 3% for Loci A and F to 23% for Locus C) for the log-linear model.

Discussion

Under the $G \times E$ simulated model presented here, it is more powerful to test for association than to test for linkage. Indeed, the MIT method has extremely poor power to detect the genetic factor either with or without taking $G \times E$ interaction into account. This could be explained by the low value of the interaction coefficient used in the simulations. Gauderman and Siegmund [5] actually showed that for an interaction coefficient less than 3 (or greater than $1/3$), the MIT will not be efficient.

For the association-based approaches, accounting for the environmental factor increases the power to detect the genetic susceptibility factor from 78% to 87% for the log-linear method and from 95% to 98% for the case-control method. This gain in power is rather limited even though under the simulated model, the gene has an effect only in exposed subjects. This could be linked to the fact that the exposure is relatively frequent in the population, as shown by Selinger-Leneman et al. [3].

If one is interested in detecting the interaction, the case-only design is shown to be the most efficient. However, its validity depends on some assumptions, in particular, the independence between both genetic and environmental factors. Type I error rates were actually higher than expected (13% instead of 5%). However, it should be noted that type I errors estimated for the two other methods were also inflated. This was essentially due to Locus C,

Table 1: Proportion of alleles shared IBD in the sib-pairs between 1500 sib pairs over 100 replicates

	π^a	π_{UU}^b	π_{EU}^c	π_{EE}^d
Average	0.502	0.500	0.501	0.503
SD ^e	0.008	0.018	0.018	0.013
Minimum	0.485	0.464	0.455	0.480
Maximum	0.525	0.543	0.543	0.539

^a π is the total proportion.

^b π_{UU} , proportion in sib pairs with 0 exposed sibs.

^c π_{EU} , proportion in sib pairs with 1 exposed sibs.

^d π_{EE} , proportion in sib pairs with 2 exposed sibs.

^eSD, standard deviation

Table 2: Power and estimates of interaction coefficients of the four tests over 100 replicates

Test	Power (%) ^a			Average interaction coefficients [95% CI]	
	G+I	G	I	I _{Bb}	I _{BB}
Mean interaction test	6	8	12	-	-
Log-linear-modeling ^b	87	78	53	1.33 [1.03–1.71]	1.72 [1.13–1.83]
Case-control ^b	98	95	69	1.39 [0.97–1.96]	1.88 [1.08–3.10]
Case-only ^b	-	-	95	1.39 [1.05–1.72]	1.86 [1.39–2.96]
Log-linear-modeling ^c	33	23	20	1.35 [0.79–2.15]	1.79 [0.82–3.68]
Case-control ^d	79	68	42	1.41 [0.85–2.09]	1.96 [0.99–3.37]

^a G+I, genetic effect accounting for interaction; G, genetic effect not accounting for interaction; I, G × E interaction.

^b Samples of 1500 families were used corresponding to 4500 (1500 triads), 3000 (1500 cases and 1500 controls), and 1500 genotyped individuals for the log-linear-modeling, the case-control and the case-only design, respectively.

^c Samples of 500 triads are considered corresponding to 1500 genotyped individuals.

^d Samples of 750 cases and 750 controls are considered here to limit the number of genotyped individuals to 1500.

which interacts with sex and might thus indirectly be associated with tobacco exposure. When Locus C is excluded, type I error rates were close to expectation with the log-linear model (5%) and with the case-control (6%), but were still increased for the case-only design (10%).

Another point of concern was the model issue. In fact, the true model was dominant but dominance was rejected in the majority of the replicates, though less often for the log-linear method than for the case-control. A plausible explanation for this distortion could be the fact that sib

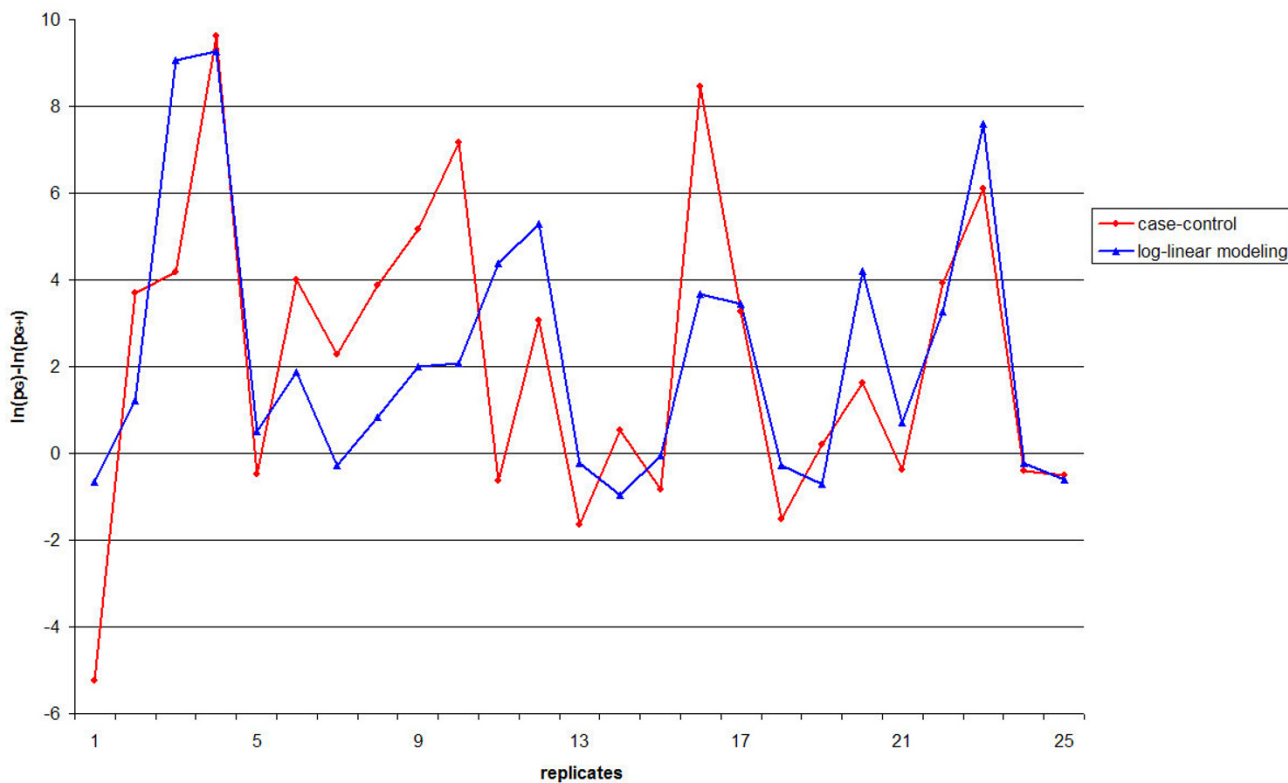


Figure 1
Difference in p-values of G+I and G tests. Difference is represented for the case-control (red plot) and the log-linear-modeling (blue plot) by $\ln(p_G) - \ln(p_{G+I})$ reported over the first 25 replicates.

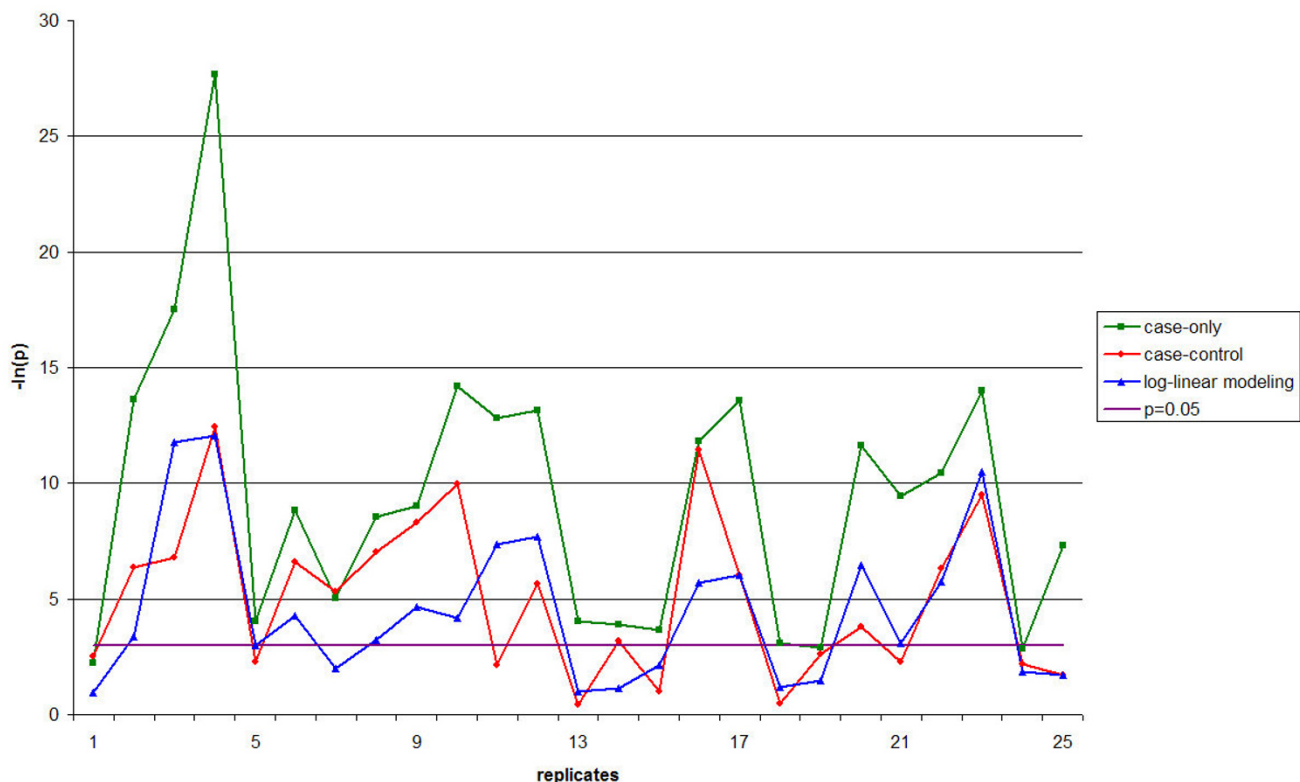


Figure 2

Comparison of the p -values of the interaction tests. $-\ln(p)$ are reported for the case-only design (green plot), the case-control design (red plot) and the log-linear-modeling method (blue plot) over the first 25 replicates.

pairs are ascertained, leading to a modification in expected parental genotype distributions. This is partially corrected for in the log-linear model by the conditioning on the parent genotypes.

All association approaches considered here do not take full advantage of the data because only one of the two sibs is used in each sib pair. It would be interesting to extend the methods to use the whole sibship while correcting for the dependence between the sibs.

Conclusion

Although this study argues in favor of the use of the case-only design to detect a $G \times E$ interaction, it shows that if one is interested in detecting gene effect, accounting for the exposure is not necessary. Of course, this depends strongly on the underlying model and could probably be linked to the high exposure frequency. It will be of interest to compare the different methods presented here using a much larger range of models.

Competing interests

The author(s) declare that they have no competing interests.

Acknowledgements

This article has been published as part of *BMC Proceedings* Volume 1 Supplement 1, 2007: Genetic Analysis Workshop 15: Gene Expression Analysis and Approaches to Detecting Multiple Functional Loci. The full contents of the supplement are available online at <http://www.biomedcentral.com/1753-6561/1?issue=S1>.

References

- Ottman R: **Gene-environment interaction: definition and study designs.** *Prev Med* 1996, **25**:764-770.
- Dizier MH, Selinger-Leneman H, Genin E: **Testing linkage and gene \times environment interaction: comparison of different affected sib-pair methods.** *Genet Epidemiol* 2003, **25**:73-79.
- Selinger-Leneman H, Genin E, Norris JM, Khat M: **Does accounting for gene-environment ($G \times E$) interaction increase the power to detect the effect of a gene in a multifactorial disease?** *Genet Epidemiol* 2003, **24**:200-207.
- Andrieu N, Goldstein AM: **Epidemiologic and genetic approaches in the study of gene-environment interaction: an overview of available methods.** *Epidemiol Rev* 1998, **20**:137-147.
- Gauderman WJ, Siegmund KD: **Gene-environment interaction and affected sib pair linkage analysis.** *Hum Hered* 2001, **52**:34-46.

6. Umbach DM, Weinberg CR: **The use of case-parent triads to study joint effects of genotype and exposure.** *Am J Hum Genet* 2000, **66**:251-261.
7. Khoury MJ, Flanders WD: **Non-traditional epidemiologic approaches in the analysis of gene-environment interaction: case-control studies with no controls!** *Am J Epidemiol* 1996, **144**:207-213.
8. Blackwelder W, Elston R: **A comparison of sib-pair linkage tests for disease susceptibility loci.** *Genet Epidemiol* 1985, **2**:85-97.

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
http://www.biomedcentral.com/info/publishing_adv.asp

