

Missense Mutations in Disease Genes: A Bayesian Approach to Evaluate Causality

Gloria M. Petersen,¹ Giovanni Parmigiani,² and Duncan Thomas³

¹Department of Epidemiology, Johns Hopkins School of Public Health, Baltimore; ²Institute of Statistics and Decision Sciences, Duke University, Durham; and ³Division of Biostatistics, Department of Preventive Medicine, University of Southern California, Los Angeles

Summary

The problem of interpreting missense mutations of disease-causing genes is an increasingly important one. Because these point mutations result in alteration of only a single amino acid of the protein product, it is often unclear whether this change alone is sufficient to cause disease. We propose a Bayesian approach that utilizes genetic information on affected relatives in families ascertained through known missense-mutation carriers. This method is useful in evaluating known disease genes for common disease phenotypes, such as breast cancer or colorectal cancer. The posterior probability that a missense mutation is disease causing is conditioned on the relationship of the relatives to the proband, the population frequency of the mutation, and the phenocopy rate of the disease. The approach is demonstrated in two cancer data sets: BRCA1 R841W and APC I1307K. In both examples, this method helps establish that these mutations are likely to be disease causing, with Bayes factors in favor of causality of 5.09 and 66.97, respectively, and posterior probabilities of .836 and .985. We also develop a simple approximation for rare alleles and consider the case of unknown penetrance and allele frequency.

Introduction

The commitment of the genetics research community to map the human genome and identify disease-causing genes has been successful in advancing our knowledge

in a number of areas, such as molecular technology leading to more-rapid cloning and sequencing of genes (Guyer and Collins 1995; Schuler et al. 1996; Savill 1997). These major advances have opened other important areas of investigation, including biochemical mechanisms for disease and genetic epidemiological implications of newly discovered genes in patients and populations (Ellsworth et al. 1997).

Of major importance is the interpretation of kinds of mutations that occur in genes and whether they may be involved in causation of disease (Cooper and Krawczak 1993). It is widely accepted that mutations that result in frameshifts (insertions or deletions) can significantly alter the protein product, as can point mutations that result in splice-site alterations or stop codons (nonsense mutations). A problematic type of point mutation is one that results in a substitution of one amino acid for another (missense mutation). These amino acid substitutions can be disease causing if they affect an important functional region of the protein. A well-known example is the substitution of valine for glutamic acid in the β -globin gene at the sixth codon, which results in an abnormal hemoglobin that is less soluble in deoxygenated blood, the basic biochemical defect in sickle-cell anemia. Alternatively, an amino acid substitution may occur in a less critical or conserved region of the protein and may be tolerated, such that it results in an isoform of the protein and, genetically, is interpreted as an allelic variant or possible polymorphism.

With the identification of disease-causing genes, a major challenge to molecular geneticists is the detection and characterization of mutations in affected persons (Cotton 1997). The definitive methodology is DNA sequencing, which will identify point mutations, including missense mutations. The interpretation of a missense mutation is often inconclusive as to whether it is disease causing. In particular, when the mutation is frequent and carriers exhibit lower penetrance, it poses a greater challenge because one could argue that it is a polymorphism.

We present a method that utilizes a Bayesian approach for families with multiple affected persons, to establish statistically whether a missense mutation in an autosomal dominant gene is disease causing. We apply this

Received December 22, 1997; accepted for publication April 7, 1998; electronically published May 15, 1998.

Address for correspondence and reprints: Dr. Gloria M. Petersen, Department of Epidemiology, Johns Hopkins University, School of Public Health, 615 North Wolfe Street, Baltimore, MD 21205. E-mail: gpetersen@jhsph.edu

© 1998 by The American Society of Human Genetics. All rights reserved. 0002-9297/98/6206-0031\$02.00

method to published data on a missense mutation of the adenomatous polyposis coli (APC) gene in familial colon cancer (Laken et al. 1997) and a missense mutation of BRCA1 in hereditary breast cancer (Barker et al. 1996).

Methods

Study Design

We propose a study design based on testing affected relatives of probands who are known to be carriers of the missense mutation. This design is efficient and provides meaningful information about the association when there is good a priori information about the rate of sporadic disease in the population and about the population frequency of the mutation. If the missense mutation is not disease causing, we expect to observe a proportion of positive affected relatives that is dictated simply by the degree of relationship between the relatives in the sample and the proband. For example, among first-degree relatives of heterozygous carrier probands of a rare mutation, we would expect about one-half to be carriers. If the mutation is disease causing, we expect to see an increased proportion, beyond one-half, depending on the penetrance and the phenocopy rate. Our approach quantifies how many more to expect and weighs the evidence in favor of causality.

In addition to its practical and ethical advantages, a design based on testing only affected relatives can be statistically efficient when reliable information about phenocopy rates is available. Swift et al. (1990) studied a simpler design in which a single affected relative per proband is tested, and developed a large-sample approximation of the confidence interval of the odds ratio of a genotype, given the disease status. They also compared the statistical efficiency of affected-only designs with that of case-control designs, in which unaffected relatives also are tested, and concluded that if the disease incidence among noncarriers is known with reasonable accuracy, testing only those affected is substantially more efficient. Although our design is more general, their overall conclusions about efficiency are likely to be applicable here.

Sampling Distributions

We consider an autosomal dominant gene and focus on a single, potentially disease-causing allele (i.e., a missense mutation in a disease-causing gene). We assume that the allele frequency in the population is known, and we denote this as “ p ,” with $1 - p = q$. We also assume that the rate of disease among noncarriers (the phenocopy rate) and the penetrance among carriers are known. We define the penetrance as $\beta = P(\text{disease}|AA) =$

$P(\text{disease}|Aa)$ and the phenocopy rate as $\varphi = P(\text{disease}|aa)$. We are interested in whether the allele is disease causing—that is, in whether carriers of the allele are at increased disease risk. The presence of a causal effect is denoted by the variable C , with $C = 1$ if the disease is caused by the mutation and $C = 2$ if it is not. So $C = 1$ corresponds to $\beta > \varphi$, whereas $C = 2$ corresponds to $\beta = \varphi$. The Appendix extends our approach to the case of unknown reduced penetrance and unknown allele frequency, with a known rate of disease for the overall population.

To assess the presence of a causal effect, we compute the Bayes factor for $C = 1$ versus $C = 2$ and the a posteriori probability that the allele is disease causing, given the observed mutation-test results from one or more pedigrees. If the penetrance, prevalence, and phenocopy rate are known, the Bayes factor is simply the ratio of the likelihoods of the observed testing results under the two hypotheses. However, by use of a Bayesian approach, extensions to unknown penetrance, prevalence, and phenocopy rate are straightforward.

For K probands, with genotypes g_{01}, \dots, g_{0K} , where g_{0k} is either AA or Aa and is fixed by design, we have corresponding n_k affected relatives who are tested for the same mutation. The genotype of relative i of proband k is g_{ik} and can be AA, aa, or Aa. The calculations are also conditioned on the relationship to the proband. To keep the notation concise, we refer to the vector of the genotypes for the relatives of proband k as $\mathbf{g}_k = (g_{1k}, \dots, g_{n_k k})$. Because of the study design, all probabilities are implicitly conditional on the proband’s genotype, on the relative being affected, and on his or her degree of relationship to the proband.

We begin by considering the contribution of a single family in the sample. To evaluate this, we need to be able to evaluate the probabilities $P(\text{observed genotypes of relatives} | C)$ for all possible combinations of observed genotypes and C . If the mutation is not disease causing ($C = 2$), the conditional probability of observed genotypes of relatives in family k is

$$\gamma_{\text{noncausal},k} \equiv P(\mathbf{g}_k | g_{0k}, C = 2) = P(\mathbf{g}_k | g_{0k}), \quad (1)$$

because affected status does not carry any information about genotype probabilities. $P(\mathbf{g}_k | g_{0k})$ depends on the allele frequency p and also on the degree of relationship to the proband, which is not explicitly incorporated into the notation but is considered in the calculation, and it can be determined simply on the basis of the degree of relationship to the proband, by use of standard conditional probability arguments (Elandt-Johnson 1971). For example, if proband k is Aa and his or her family contains only one other affected relative, a sibling also with genotype Aa, then $n_k = 1$ and

$$P(g_{1k}|g_{0k}, C = 2) = P(Aa|Aa, C = 2) = \frac{1}{2}(1 + pq) .$$

For a heterozygous parent of the proband, we would have

$$P(g_{1k}|g_{0k}, C = 2) = P(Aa|Aa, C = 2) = \frac{1}{2} ;$$

for a homozygous sib,

$$P(g_{1k}|g_{0k}, C = 2) = P(AA|Aa, C = 2) = \frac{1}{4}p(1 + p) ;$$

and so forth. A general algorithm is presented by Li and Sacks (1954). The software package LINKAGE (Genetic Linkage Analysis, Laboratory of Statistical Genetics, Rockefeller University) can be used to automate these calculations.

If the mutation is disease causing ($C = 1$), the expected fraction of carriers among affected relatives is higher than if it is not disease causing. Conditional on a genotype vector g_k with n_k^{aa} relatives with aa genotype and under the assumption that the disease outcomes of relatives are independent given their genotypes, the probability that all relatives have the disease is $\varphi^{n_k^{aa}}\beta^{n_k - n_k^{aa}}$. By use of the genotype probabilities $P(g_k|g_{0k})$ as priors, we can determine genotype probabilities under $C = 1$ via Bayes rule. For the $C = 1$ case, the probability of the observed genotypes of relatives in family k is

$$\begin{aligned} \gamma_{causal,k} &= P(g|g_{0k}, C = 1) \\ &= \frac{P(g_k|g_{0k})\varphi^{n_k^{aa}}\beta^{n_k - n_k^{aa}}}{\sum_g P(g|g_{0k})\varphi^{n_g^{aa}}\beta^{n_k - n_g^{aa}}} \end{aligned} \quad (2)$$

where \sum_g ranges over all 3^k possible genotype combinations and n_g^{aa} is the number of aa elements in the vector g . In practice, the number of terms in the summation in the denominator can be reduced by factoring terms that lead to the same value of n_g^{aa} . For example, for a proband with one sibling who is Aa, we have

$$\begin{aligned} P(g_{1k}|g_{0k}, C = 1) &= P(Aa|Aa, C = 1) \\ &= \frac{P(Aa|Aa)P(disease|Aa \text{ or } AA)}{P(Aa \text{ or } AA|Aa)P(disease|Aa \text{ or } AA) + P(aa|Aa)P(disease|aa)} \\ &= \frac{(1 + pq)/2}{1 - q(1 + q)/4 + q(1 + q)\varphi/4} . \end{aligned}$$

When the disease has a variable age at onset and when the age-at-onset data for relatives are available, β and φ should be specified in terms of age-specific survival contributions. This would entail taking each subject's age into account. However, by rewriting expression (2) as

$$\gamma_{causal,k} = \frac{P(g_k|g_{0k})(\varphi/\beta)^{n_k^{aa}}}{\sum_g P(g|g_{0k})(\varphi/\beta)^{n_g^{aa}}} ,$$

it can be seen that $\gamma_{causal,k}$ depends on β and φ only via their ratio and therefore is valid as long as the ratio of the penetrance to the phenocopy rate remains constant with respect to age. It is possible to carry out an age-dependent analysis along the lines of the approach described here. However, this requires either a large number of observations or good a priori knowledge of the penetrance function for the mutation under consideration. These were not available for the applications considered in this article.

Probability of a Disease-Causing Effect

We can now evaluate the probability of a causal effect conditional on the observed mutation tests in the affected relatives ("data"). Because the families can be considered independent, in our case

$$\begin{aligned} P(\text{data}|C = 1) &= \prod_{k=1}^K \gamma_{causal,k} \quad \text{and} \\ P(\text{data}|C = 2) &= \prod_{k=1}^K \gamma_{noncausal,k} , \end{aligned}$$

so that the Bayes factor in favor of the hypothesis of causality is

$$B = \frac{\prod_{k=1}^K \gamma_{causal,k}}{\prod_{k=1}^K \gamma_{noncausal,k}} .$$

Using expressions (1) and (2) and algebraic manipulation, we can express the Bayes factor simply in terms of the genotype coefficients and the penetrance-to-phenocopy rate ratio:

$$1/B = \prod_{k=1}^K \sum_g P(g|g_{0k}) \left(\frac{\varphi}{\beta}\right)^{n_g^{aa} - n_k^{aa}} .$$

This incorporates the dependence among relatives of the same proband and the possibility that there is more than one copy of the mutated allele in the same family.

By using Bayes rule, which states

$$\begin{aligned} P(C = 1|\text{data}) &= \frac{P(C = 1)P(\text{data}|C = 1)}{P(C = 1)P(\text{data}|C = 1) + P(C = 2)P(\text{data}|C = 2)} , \end{aligned}$$

and by denoting the a priori odds in favor of causality by $O = P(C = 1)/P(C = 2)$, the posterior probability of causality is $P(C = 1|\text{data}) = OB/(OB + 1)$.

The choice of a priori odds can be important and must

Table 1
Data for the Colon Cancer Example

Proband	Affected Relative(s)	Genotype(s) of Relative(s)		$\gamma_{causal,k} / \gamma_{noncausal,k}$		
		Relative(s)	$\gamma_{causal,k}$	$\gamma_{noncausal,k}$		
1	Two siblings and one niece from a third sibling	Aa, Aa, Aa	.42	.08	5.02	
2	Sibling	Aa	.81	.52	1.56	
3	Mother	Aa	.79	.50	1.58	
4	Sibling	aa	.18	.47	.38	
5	Sibling	Aa	.81	.52	1.56	
6	Grandmother	Aa	.61	.28	2.15	
7	Mother	Aa	.79	.50	1.58	
8	One sibling and her offspring	Aa, Aa	.70	.26	2.72	

NOTE.—Summary of genetic-testing results and contributions of each family to the likelihood ratio. The rightmost column represents the contribution of each family to the calculation of the Bayes factor, which is the product of the values in the column.

reflect the context of the analysis. If the mutation was selected because it resides on a known disease-causing gene, the prior odds may be high, while, when screening for a disease-causing mutation, a low population-based prior probability of causality may be more appropriate. In any case, analysis of the Bayes factor gives a measure of the weight of evidence of the data in favor of the hypothesis that does not require specification of prior odds. Kass and Raftery (1995) give an in-depth discussion of interpretation and calibration of Bayes factors. A convenient option is the assumption of $P(C = 1) = \frac{1}{2}$, which gives $O = 1$ —that is, even a priori odds to the hypothesis of causality.

Results

Familial Breast Cancer and BRCA1 R841W

We investigate the causality of the BRCA1 R841W mutation, using two pedigrees from the study by Barker et al. (1996). The disease phenotype of interest is either breast or ovarian cancer. There are $K = 2$ probands who have tested, affected relatives. Proband 1160 has one heterozygous sibling affected with breast cancer, $g_1 = g_{11} = Aa$. Proband 728 has two heterozygous siblings affected with breast cancer. Thus, $g_2 = (g_{12}, g_{22}) = (Aa, Aa)$. We assume a phenocopy rate $\varphi = .125$, a penetrance $\beta = .85$, indicated as plausible for other BRCA1 mutations (Ford and Easton 1995), and an allele frequency $p = 1/100$. Evidence for the estimation of age-specific penetrance for this specific mutation was insufficient. Our calculations are based on the assumption that the ratio β/φ of penetrance, for carriers and non-

carriers, does not depend on age. We will later assess sensitivity to p and β/φ . Not all affected family members were tested; we assume that the reason why they were not tested is unrelated to their genotype.

By use of expressions (1) and (2),

$$\gamma_{noncausal,1} = \frac{1}{2}(1 + pq) = .505 ;$$

$$\gamma_{noncausal,2} = \frac{1}{4}(5pq + p^2 + q^2) = .257 ;$$

$$\gamma_{causal,1} = \gamma_{noncausal,1} / [1 - \frac{1}{4}q(1 + q)(1 - \varphi)] = .871 ; \text{ and}$$

$$\begin{aligned} \gamma_{causal,2} = & \gamma_{noncausal,2} / \{ \frac{1}{8}q(1 + q)\varphi^2 + p^2 + \frac{1}{4}q^2 + \frac{25}{16}pq \\ & + [1 - \frac{1}{8}q(1 + q) - p^2 - \frac{1}{4}q^2 - \frac{25}{16}pq]\varphi \} \\ = & .761 . \end{aligned}$$

The evidence from family 1 is $\sim 1.73 \times$ more likely under the hypothesis of causality. The evidence from family 2 is $2.96 \times$ more likely under the hypothesis of causality. This results in a Bayes factor of 5.09, indicating that test results from the two families are $\sim 5 \times$ more likely under the hypothesis of causality than under the hypothesis of no causality. If prior odds for causality are even, the posterior probability of causality is .836.

Our results are virtually unchanged as p varies over the range 0–.01 and decreases to .82 if p is set to .05. The rare-allele approximation discussed in the Appendix gives virtually the same results. The results are only mildly sensitive to assumptions about penetrance. At a penetrance of 1, the probability of causality is .85. At a penetrance of .5, it is .8.

Familial Colon Cancer and APC I1307K

Mutations of the APC gene can be associated with increased risk of colorectal cancer and colorectal adenomas. We investigate the causality of mutation T→A at APC nucleotide 3920, using eight pedigrees reported by Laken et al. (1997). We used a disease rate α of .2, on the basis of clinical judgment, and an allele frequency p of the mutation of .036, on the basis of table 1 in the report by Laken et al. (1997); we estimated allele frequencies separately in the disease group and control group, and combined the two estimates using the postulated value of α . Expressions (1) and (2) were computed for each proband. The results are summarized in table 1.

The resulting Bayes factor is 66.97, which is strong evidence in favor of causality. The posterior probability under even prior odds is .985. The rare-allele approximation discussed in the Appendix produces a Bayes factor of 194.90, leading to a posterior probability of .995.

Because the mutation is relatively common, the rare-allele approximation does not work as well as in the breast cancer example. Also, use of expression (A3) leads to ignoring the dependence between proband 8's sibling and that sibling's offspring and therefore is not appropriate in this case.

Following the development discussed in the first section of the Appendix, we also conducted an analysis with unknown penetrance and allele frequency. We assumed that causality ($C = 1$) corresponds to $\beta > \varphi$, and we assumed that, if the mutation is indeed causal, then all values of β are equally likely a priori. With regard to p , we used information from table 1 in the report by Laken et al. (1997) to specify a prior on the fraction of carriers and then converted that into a distribution for p . Because the fraction of diseased individuals in the sample approximates the overall population fraction, we combined diseased and nondiseased individuals when specifying our prior. The resulting specification is $p = 1 - (1 - r)^{1/2}$, where r has a beta distribution with parameters $47 + 22 = 69$ and $766 + 211 - 47 - 22 = 908$.

To analyze the effect of uncertainty about β and p on the Bayes factor and posterior probability of $C = 1$, we first computed the a posteriori probability distribution $p(\beta, p | \text{data}, C = 1)$. This inference is conditional on the assumption of uninformative ascertainment, which may be violated in this data set. Although inference of the penetrance and prevalence is not the focus of the methodology proposed here, incorporation of uncertainty via $p(\beta, p | \text{data}, C = 1)$ is an effective strategy to make inferences about causality without relying on strong assumptions of β and p . To illustrate the range of plausible values of β , figure 1 graphs the a posteriori probability distribution of β conditional on $p = .035$. For the purpose of testing for causality, the important aspect of this distribution is that it assigns very low probability to values close to the phenocopy rate and high probability to values far away from it, so that, even though the actual magnitude of β is not accurately estimated, reliable conclusions about causality can be reached. We computed the Bayes factor and posterior probability of $C = 1$ using expression (A2). The Bayes factor is 53.21, and the resulting posterior probability is .982. Even in the presence of substantial uncertainty about the exact value of the prevalence parameter, the evidence in favor of causality in this example remains strong. Results are somewhat sensitive to the specification of the phenocopy rate α , but the evidence in favor of causality is not questioned within a broad range of values. At $\alpha = .1$, the posterior probability is .987, whereas at $\alpha = .3$ it is .965.

Discussion

We propose an approach, using Bayesian methods, to statistically determine whether a missense mutation in

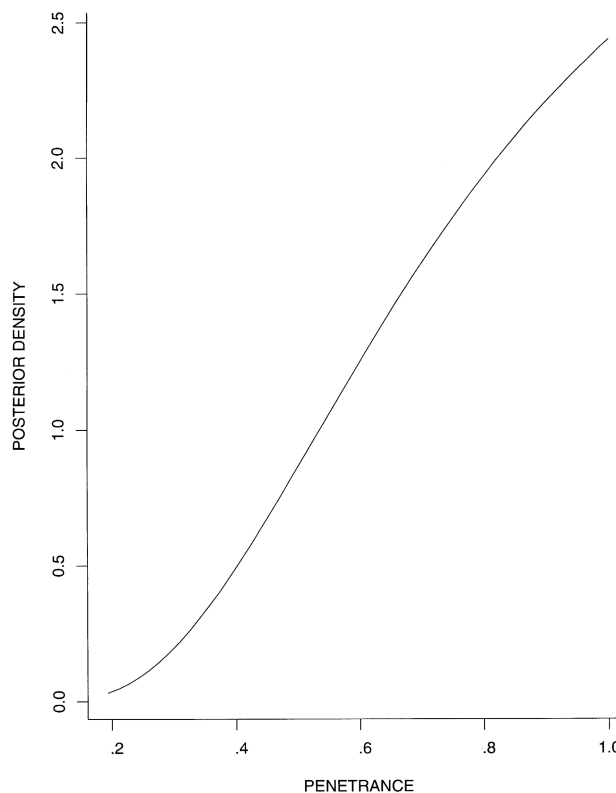


Figure 1 Posterior distribution of penetrance parameter β

an autosomal dominant disease gene is likely to be disease causing. Our method is particularly useful in evaluating common disease phenotypes, only a small proportion of which may be caused by known disease-causing genes, such as those for hereditary breast cancer or familial colorectal cancer. This approach employs the testing of other affected relatives for the missense mutation in pedigrees that have been identified through a proband with the missense mutation. The posterior probability that the mutation is disease causing is conditioned on the relationship of the relative to the proband, the frequency of the mutation, and the phenocopy rate of the disease. Because only affected relatives are selected for genetic analysis, this method is efficient, from the perspective of conducting a study. Because unaffected relatives are not studied, this method avoids the pitfalls that relate to recruitment of such persons for research in the current genetic-testing environment (Hubbard and Lewontin 1996; Geller et al. 1997; Holtzman et al. 1997).

However, by limiting analysis only to families with at least two affected relatives, there is a potential for biasing the conclusion in favor of causality, since the penetrance may, in reality, be lower than the sample indicates. This can be partially addressed by varying penetrance esti-

mates in the computations. This design element also may affect our results if there is a sizable proportion of families that have only one affected person (the proband). If the unaffected phenotypes in these families are caused by reduced penetrance, the magnitude of the effect for specific mutations may be empirically assessed only by genotyping unaffected relatives. This design modification, of course, engenders logistical issues for the execution of the study.

Another important consideration is the potential bias that may occur in the selection of affected relatives for genotyping. Our approach assumes that inclusion in the study is unrelated to genotype. Ideally, any and all affected biological relatives should be genotyped, but availability of DNA can be affected by numerous factors (relatives who are deceased or who are unable or unwilling to participate). If availability for study is non-random, for whatever reasons, this may have an effect on the posterior probability, but the effect would depend on the nature of the selection bias.

We have shown that this method can be applied usefully to empirical data. We analyzed families identified through probands with the APC I1307K mutation (Laken et al. 1997), in which 10 of 11 affected relatives tested positive for the same mutation, including three second-degree relatives. We found a .985 posterior probability that the APC I1307K mutation is causally related to the colorectal neoplasia in these families. In the case of BRCA1 R841W (Barker et al. 1996), there were fewer families and persons studied, but we estimated a .836 posterior probability that this missense mutation may be causally related to breast cancer or ovarian cancer in these families.

These two examples contrast another challenge in the interpretation of missense mutations—namely, a plausible mechanism of disease causation engendered by the mutation (Ellsworth et al. 1997; Fearon 1997). In the case of APC I1307K, the T→A transversion in the APC gene results in an (A)⁸ tract, which appears to engender an inherent instability in the gene, allowing deleterious mutations to occur in this critical gene during subsequent cell division (Laken et al. 1997). In the case of BRCA1 R841W, however, it is difficult to posit a plausible mechanism because the function of the gene remains unclear. Our proposed method may provide additional evidence to support a disease-causing missense mutation in such cases. That a missense mutation may have a plausible mechanism for causing disease can justify a more realistic estimate of the penetrance (whether higher or lower), whereas a missense mutation in a gene of unknown function may conservatively be justified to be lower.

There is a possibility, although remote, that the missense mutation may be in linkage disequilibrium with a “true” deleterious mutation elsewhere in the allele. Con-

sidering the way in which many missense mutations are identified through DNA sequencing, such a possibility would have been uncovered. In the case of our colon cancer example, Laken et al. (1997) reported that the APC I1307K-bearing alleles in two carriers were fully sequenced and that no other mutations were detected.

This approach is predicated on the assumption that a single missense mutation has been singled out for particular attention, on a priori grounds. Since disease genes may potentially be highly polymorphic, additional complications arise for the interpretation of the causality of polymorphisms that may come to attention because of their frequency in multiple case families, since such polymorphisms are more likely than randomly selected polymorphisms to appear to be causal, even if $C = 2$ were true. A hierarchical Bayesian treatment of this problem might entail simultaneous consideration of all known polymorphisms, with additional parameters for variation in β and p between polymorphisms. The model could include “prior covariates,” such as the position in the gene or the nature of the particular amino acid substitution that results, for which the effects on penetrance are to be estimated empirically. Such an analysis is beyond the scope of this article, but similar methods have been applied in the treatment of multiple-exposure problems in epidemiology (Greenland 1993) and gene-environment interactions (Aragaki et al. 1997).

The design approach of the study of affected relatives of a proband was employed by Swift et al. (1990) and was applied by Athma et al. (1996) in a study of breast cancer risk in ataxia telangiectasia heterozygotes. This approach differs from the method proposed here, in a number of elements. Their design requires ascertainment of one affected relative of an index carrier (who may not be affected with the disease of interest) per family. Inference and testing are based on a large-sample approximation of the confidence interval on the odds ratio of the genotype, given the disease status. An important advantage of the Bayesian approach is that it is simple to derive inferences, without having to rely on large-sample approximations that could be inaccurate, and it is straightforward to incorporate uncertainty about unknown additional parameters, as discussed in the Appendix.

Finally, we have developed a simple program in S-PLUS, shown in the Appendix, which provides computation of the Bayes factor and posterior probabilities in the rare-allele case. This basic approach can be extended to consider reduced penetrance with variable age at onset and errors in genetic testing. Development of these extensions is currently in progress.

Conclusions

The approach proposed here anticipates one of the outgrowths of the Human Genome Project and genetic-

disease research: a number of families that segregate common disease phenotypes may be identified to carry missense mutations of known disease-causing genes. A concomitant problem in interpretation of the significance of the missense mutations is determination of whether they are disease causing or simple polymorphisms. We have developed a logistically efficient and computationally feasible approach that may more quickly help determine the importance of such mutations.

Acknowledgments

The authors thank Stephen Gruber, Steve Laken, Bert Vogelstein, Kenneth Kinzler, Frank Giardiello, Stan Hamilton, and Susan Booker for their contributions to this study. This research was supported, in part, by NIH National Cancer Institute grants CA 52862, R01 CA 63721, and P50 CA 62924 (Specialized Program of Research Excellence in Gastrointestinal Cancer, Johns Hopkins University) and P50 CA68438 (Specialized Program of Research Excellence in Breast Cancer, Duke University) and by the Clayton Fund. This work was carried out while Dr. Parmigiani was visiting the Department of Biostatistics, Johns Hopkins University; the department’s warm hospitality is gratefully acknowledged.

Appendix

Statistical Details

Unknown Penetrance and Prevalence

The derivation of the posterior probability of a causal effect can be extended to the case of unknown penetrance and prevalence parameters. Within a Bayesian approach, this is carried out by replacing the fixed values of β and p with prior probability distributions reflecting information from prior studies or other biological evidence. The prior distribution of the unknown β and p is denoted here by $\pi(\beta, p)$. If the penetrance of a similar mutation has been previously studied, published data can be used to specify an informative prior for β . Otherwise, an attractive option is to assume a priori that all values $>\varphi$ are equally likely, leading to $\pi(\beta) = 1/(1 - \varphi)$, $\varphi < \beta \leq 1$. Because the mutation may be responsible for a nonnegligible fraction of the overall cases, the assumption that φ is known also needs to be relaxed. Here, we assume that the overall incidence α is known. The overall incidence depends on φ , β , and p , via the relationship $\alpha = q^2\varphi + (1 - q^2)\beta$, leading to $\varphi = [\alpha - (1 - q^2)\beta]/q^2$.

As in the section Sampling Distributions, under $C = 1$ it is natural to assume that $\beta > \varphi$. With this specification, then, we are testing a point null hypothesis

against a composite alternative hypothesis. Berger and Delampady (1987) review the Bayesian approach to this problem and compare it with frequentist approaches. For a fixed α , a prior distribution of β and p implies a prior distribution of φ . Also, $\varphi < \beta$ whenever $\alpha < \beta$.

Under these assumptions, expression (1) remains the same, whereas expression (2) can now be interpreted as conditional on β and p ; that is,

$$\begin{aligned} \gamma_{\text{causal},k}(\beta, p) &= P(\mathbf{g}|\mathbf{g}_{0k}, C = 1, \beta, p) \\ &= \frac{P(\mathbf{g}_k|\mathbf{g}_{0k})\beta^{n_k - n_k^{\text{aa}}}\varphi^{n_k^{\text{aa}}}}{\sum_{\mathbf{g}} P(\mathbf{g}|\mathbf{g}_{0k})\beta^{n_k - n_k^{\text{aa}}}\varphi^{n_k^{\text{aa}}}}; \end{aligned} \quad (\text{A1})$$

rewriting φ in terms of α , β , and p ,

$$\gamma_{\text{causal},k}(\beta, p) = \frac{P(\mathbf{g}_k|\mathbf{g}_{0k})\{[\alpha - (1 - q^2)\beta]/q^2\}^{n_k^{\text{aa}}}}{\sum_{\mathbf{g}} P(\mathbf{g}|\mathbf{g}_{0k})\{[\alpha - (1 - q^2)\beta]/q^2\}^{n_k^{\text{aa}}}},$$

so that $P(\text{data}|C = 1) = \prod_{k=1}^K \gamma_{\text{causal},k}(\beta, p)$.

The Bayes factor for the hypothesis of causality is now

$$B = \frac{\int_0^1 \int_{\alpha}^1 \prod_{k=1}^K \gamma_{\text{causal},k}(\beta, p) \pi(\beta, p) d\beta dp}{\prod_{k=1}^K \gamma_{\text{noncausal},k}}, \quad (\text{A2})$$

and the posterior probability incorporating uncertainty about β and p can be derived in the usual way as $P(C = 1|\text{data}) = OB/(1 + OB)$, where O is the prior odds in favor of $C = 1$. Expression (A2) is easily evaluated numerically by use of a Monte Carlo approach. To evaluate the numerator, it is sufficient to generate a sample for the prior distribution $\pi(\beta, p)$ and to compute the average of the resulting values of $\prod_{k=1}^K \gamma_{\text{causal},k}(\beta, p)$.

Expression (A1) can be used to derive an a posteriori distribution on β and p , conditional on $C = 1$:

$$p(\beta, p|\text{data}, C = 1) = \frac{\prod_{k=1}^K \gamma_{\text{causal},k}(\beta, p) \pi(\beta, p)}{\int_0^1 \int_{\alpha}^1 \prod_{k=1}^K \gamma_{\text{causal},k}(\beta, p) \pi(\beta, p) d\beta dp}.$$

This provides a valid inference on the penetrance if it can be assumed that the selection mechanism—that is, choosing families with at least one affected relative and eliminating index cases—is not inducing a bias in the estimation of penetrance.

Rare Alleles

When the allele is rare and there are no homozygous carriers in the sample, ignoring the dependence among relatives of the same proband and ignoring the possibility that there is more than one copy of the mutated allele in the same family have a limited impact on the final answer. Expressions (1) and (2) then can be sim-

plified considerably. It is interesting to study the form of these simplified expressions. All the $N = \sum n_k$ relatives can be considered to be approximately independent. For relative j , let r_j be the degree of relationship to the proband and b_j be an indicator variable of whether or not relative j is a carrier. Then

$$P(\text{data}|C = 2) \approx \prod_{j=1}^N \left(\frac{1}{2}\right)^{r_j b_j} [1 - \left(\frac{1}{2}\right)^{r_j}]^{1-b_j}$$

$$= \prod_{j=1}^N \frac{(2^{r_j} - 1)^{1-b_j}}{2^{r_j}} .$$

Applying Bayes rule for each individual and using independence, we also obtain

$$P(\text{data}|C = 1) \approx \prod_{j=1}^N \frac{(1/2)^{r_j b_j} \beta^{b_j} \{ [1 - (1/2)^{r_j}] \varphi \}^{1-b_j}}{(1/2)^{r_j} \beta + [1 - (1/2)^{r_j}] \varphi}$$

$$= \prod_{j=1}^N \frac{\beta^{b_j} [2^{r_j} - 1] \varphi^{1-b_j}}{1 + (2^{r_j} - 1) \varphi} .$$

After some simple manipulations, these lead to a Bayes factor in favor of a causal effect of

$$B = \prod_{j=1}^N \frac{(\beta/\varphi)^{b_j} 2^{r_j}}{(\beta/\varphi) + 2^{r_j} - 1} . \tag{A3}$$

Each of the N terms is the ratio of the probability of the evidence under $C = 2$ to the probability of the evidence under $C = 1$, or the weight of evidence against $C = 2$. For example, for a first-degree relative, the ratio is $2\beta/(\beta + \varphi)$, which is >1 , if the relative is a carrier; and the ratio is $2\varphi/(\beta + \varphi)$, which is <1 , if the relative is not a carrier.

In the unknown penetrance case, the expression for the Bayes factor becomes

$$B = \int_{\alpha}^1 \prod_{j=1}^N \frac{(\beta/\varphi)^{b_j} 2^{r_j}}{(\beta/\varphi) + 2^{r_j} - 1} \pi(\beta) d\beta .$$

A function evaluating this expression by use of a Monte Carlo algorithm is provided in figure A1. In the rare-allele case, it is simple to develop general-purpose functions for computing the Bayes factors and the posterior probability of causality. Here we present a function, written in the statistical package S-PLUS, that handles both the known and unknown penetrance cases (this function is available at <http://www.isds.duke.edu/~gp>). The inputs to the function are the following: the vector RELATIONS, with as many elements as there are relatives and with each element a 1, 2, etc., for first-degree

```
arcs.rare _ function(relations,genotypes,phenocopy.rate,
                    prevalence=1,unknown.prevalence=F,
                    MonteCarlo=1000) {

  rr _ relations
  gg _ genotypes
  phi _ phenocopy.rate
  beta _ prevalence
  M _ MonteCarlo
  NN _ length(relations)
  integrand _ NULL

  if(!unknown.prevalence) {
    bf _ exp ( sum(rr*log(2) +
                  gg*log(beta/phi) -
                  log( beta/phi + 2^rr - 1 ) ) )
  }

  if(unknown.prevalence) {
    beta _ runif(M,phi,1)
    for (m in 1:M) {
      integrand[m] _ exp ( sum (rr*log(2) +
                                gg*log(beta[m]/phi) -
                                log( beta[m]/phi + 2^rr - 1 ) ) )
    }
    bf _ mean(integrand)
  }

  post _ bf / ( bf + 1 )
  return(bf,post)
}
```

Figure A1 S-PLUS function for computing the Bayes factor and posterior probability of a disease-causing effect, for the rare-allele case. This function is available at <http://www.isds.duke.edu/~gp>.

relatives, second-degree relatives, etc.; the vector GENOTYPES, again with one element for each relative, either 0 if the relative is aa or 1 if the relative is Aa (use of the rare-allele approximation is not recommended when there are homozygous relatives in the sample); the scalar PREVALENCE, for the value of the prevalence, if known; the Boolean variable UNKNOWN.PREVALENCE, which can be set to T or F, depending on whether the penetrance is known (if UNKNOWN.PREVALENCE is set to T, the input value of PREVALENCE is ignored); and the integer MONTECARLO, specifying the number of Monte Carlo samples desired in the evaluation of the Bayes factor when the penetrance is unknown.

Electronic-Database Information

URLs for data in this article are as follows:

Genetic linkage analysis, Laboratory of Statistical Genetics, Rockefeller University, New York, <http://linkage.rockefeller.edu> (for LINKAGE)

References

Aragaki C, Greenland S, Probst-Hensch N, Haile RW (1997) Hierarchical modeling of gene-environment interactions: estimating NAT2* genotype-specific dietary effects on ade-

- nomatous polyps. *Cancer Epidemiol Biomarkers Prev* 6: 307–314
- Athma P, Rappaport R, Swift M (1996) Molecular genotyping shows that ataxia-telangiectasia heterozygotes are predisposed to breast cancer. *Cancer Genet Cytogenet* 92:130–134
- Barker DE, Almeida ERA, Casey G, Fain PR, Liao SY, Masunaka I, Noble B, et al (1996) BRCA1 R841W: a strong candidate for a common mutation with moderate phenotype. *Genet Epidemiol* 13:595–604
- Berger JO, Delampady M (1987) Testing precise hypotheses. *Stat Sci* 2:317–335
- Cooper DN, Krawczak M (1993) Human gene mutation. BIOS Scientific Publishers, Oxford
- Cotton RGH (1997) Mutation detection. Oxford University Press, New York
- Elandt-Johnson RC (1971) Probability models and statistical methods in genetics. Wiley & Sons, New York
- Ellsworth DL, Hallman DM, Boerwinkle G (1997) Impact of the Human Genome Project on epidemiologic research. *Epidemiol Rev* 19:3–13
- Fearon ER (1997) Human cancer syndromes: clues to the origin and nature of cancer. *Science* 278:1043–1050
- Ford D, Easton DF (1995) The genetics of breast and ovarian cancer. *Br J Cancer* 72:805–812
- Geller G, Botkin JR, Green MJ, Press N, Biesecker BB, Wilfond B, Grana G, et al (1997) Genetic testing for susceptibility to adult-onset cancer: the process and content of informed consent. *JAMA* 277:1467–1474
- Greenland S (1993) Methods for epidemiologic analysis of multiple exposures: a review and comparative study of maximum-likelihood, preliminary-testing, and empirical-Bayes regression. *Stat Med* 12:717–736
- Guyer MS, Collins FS (1995) How is the Human Genome Project doing, and what have we learned so far? *Proc Natl Acad Sci USA* 92:10841–10848
- Holtzman NA, Murphy PD, Watson MS, Barr PA (1997) Predictive genetic testing: from basic research to clinical practice. *Science* 278:602–605
- Hubbard R, Lewontin RC (1996) Pitfalls of genetic testing. *New Engl J Med* 334:1192–1194
- Kass RE, Raftery AE (1995) Bayes factors. *J Am Stat Assoc* 90:773–795
- Laken SJ, Petersen GM, Gruber SB, Oddoux C, Ostrer H, Giardiello FM, Hamilton SR, et al (1997) Familial colorectal cancer in Ashkenazim due to a hypermutable tract in APC. *Nat Genet* 17:79–83
- Li CC, Sacks L (1954) The derivation of joint distribution and correlation between relatives by the use of stochastic matrices. *Biometrics* 10:347–360
- Savill J (1997) Molecular genetic approaches to understanding disease. *BMJ* 314:126–129
- Schuler GD, Boguski MS, Stewart EA, Stein LD, Gyapay G, Rice K, White RE, et al (1996) A gene map of the human genome. *Science* 274:540–546
- Swift M, Kupper LL, Chase CL (1990) Effective testing of gene-disease associations. *Am J Hum Genet* 47:266–274