1 **Date:** 16 Feb 2023

2 **Title:** The effect of ascertainment on penetrance estimates for rare variants: implications for

3 establishing pathogenicity and for genetic counselling

4

5 **Running title:** Ascertainment effects on penetrance estimates

6

7 **Authors:** Andrew D. Paterson*[1,2], Sang-Cheol Seok[3], Veronica J. Vieland[3,4]

8 **ORCID:**

9 ADP: 0000-0002-9169-118X

10 SCS: 0000-0003-3955-9891

11 VJV: 0000-0002-3004-3840

12 **Affiliations:**

13    1. Program in Genetics and Genome Biology, The Hospital for Sick Children, Toronto,

14       Ontario, M5G 1X8, Canada

15    2. Divisions of Epidemiology and Biostatistics, Dalla Lana School of Public Health,

16       University of Toronto, Toronto, Ontario, Canada

17    3. Mathematical Medicine LLC, Chicago, IL

18    4. Departments of Pediatrics and Biostatistics (Emerita), The Ohio State University,

19       Columbus, OH

20 **\*Correspondence:** andrew.paterson@sickkids.ca

21 **Keywords:**

22 Penetrance, ascertainment, Mendelian disease, multiplex family, rare variant, variant

23 pathogenicity

24    **Abstract:**

25    Next-generation sequencing has led to an explosion of genetic findings for many rare diseases.

26    However, most of the variants identified are very rare and were identified in small pedigrees,

27    which creates challenges in terms of penetrance estimation and translation into genetic

28    counselling in the setting of cascade testing. We use simulations to show that for a rare

29    (dominant) disorder where a variant is identified in a small number of small pedigrees, the

30    penetrance estimate can both have large uncertainty and be drastically inflated, due to underlying

31    ascertainment bias. We have developed PenEst, an app that allows users to investigate the

32    phenomenon across ranges of parameter settings. We also illustrate robust ascertainment

33    corrections via the LOD score, and recommend a LOD-based approach to assessing

34    pathogenicity of rare variants in the presence of reduced penetrance.

35 Next-generation sequencing has led to an explosion in the number of genetic findings for many

36 rare diseases. For certain types of rare coding variants (e.g. missense, or protein truncating), if

37 the variant is sufficiently rare and has bioinformatic predictions that are severe, current

38 algorithms result in it being classified as pathogenic (1). However, the analysis of large-scale

39 sequencing from cohorts, such as ExAC (2), gnomAD (3), and the UK Biobank (4), has shown

40 that many such variants may often lack clinically significant impact. For example, ExAC

41 estimated that individuals from population cohorts carried a mean of 53 variants previously

42 thought to be sufficient causes of Mendelian diseases. Additionally, 88% of such variants had

43 MAF>1%, implying that they are likely not sufficient causes. This may indicate that such

44 variants are not causally related to disease, or perhaps, that they are causally related but with

45 reduced penetrance.

46

47 Penetrance plays an important role in understanding disease pathology, in the appropriate

48 classification of pathogenic variants, and perhaps above all in the context of genetic counseling.

49 However, most of the variants reported to date have been very rare and identified in small sets of

50 unrelated individuals (sometimes just one) or small pedigrees. Penetrance cannot be estimated

51 from a single case, or a single parent-offspring trio presenting with a *de novo* mutation in the

52 offspring. But even with multiple cases or families, determination of the penetrance can present

53 challenges. Here we focus on one such challenge: ascertainment.

54

55 Typically a variant of interest is first identified in one individual with a given phenotype.

56 Investigators may then sequence either additional relatives of the individual, or additional

57 individuals or families presenting with the same or closely related phenotypes, with the goal of

58    bolstering the case for pathogenicity. Thus, ascertainment of individuals to be sequenced

59    typically proceeds in stages. The precise ascertainment process used to enrol individuals and/or

60    families is usually at least to some extent unsystematic, and may vary between families.

61    Ascertainment is therefore challenging to model when attempting to estimate the penetrance of a

62    variant.

63

64    One situation in which ascertainment can be easily handled is "single" ascertainment, in which

65    the probability of an affected individual being ascertained is proportional to the number of

66    affected individuals in the family (5). In fact, much of the literature on inferring pathogenicity or

67    estimating penetrance tends to assume single ascertainment, e.g., (6), where ascertainment is

68    addressed by conditioning on "the proband," a procedure which is strictly correct only under true

69    single ascertainment. While it is true that the typical study ascertains families through one

70    individual who may be designated as the single "proband", this does not ensure that the study

71    meets the proportionality requirement of single ascertainment. This requirement would be

72    violated, e.g., if families with four affected members were more than twice as likely to be

73    recruited as families with just two; or, if the probability of a second sibling being ascertained

74    were dependent on the ascertainment status of the first. And in general, if either (i) ascertainment

75    is not truly single, or (ii) even if it is, if an appropriate ascertainment correction is not

76    incorporated into the estimation method, then penetrance estimates will be biased. Here we

77    consider the magnitude of that bias, across a range of plausible ascertainment models and

78    varying amounts of available data.

79

80    We focus here on sibship data. The impact of ascertainment for more complex pedigrees can be

81    approximated by considering large sibship sizes. For simplicity, we assume all parents are

82   phenotypically and genotypically unknown; including parental information does not

83   substantively affect results. We assume a very rare variant of interest (VOI), and an autosomal

84   dominant disease D. Let a qualifying individual (QI) be anyone who is both heterozygous (HET)

85   for the VOI and also affected (AFF) with D. Let $r$ be the number of QI sibs within a family, and

86   let $t$ be the number of AFF sibs regardless of VOI genotype. We also assume that, regardless of

87   VOI status, an individual might develop D due to other factors, which might be genetic

88   (involving one or more VOIs at other loci or other variants within the same gene) and/or

89   environmental (e.g., due to infections). Let $\gamma$ be the combined penetrance across all causes other

90   than the VOI under study. Since we assume the VOI is very rare, $\gamma$ is effectively the population

91   prevalence of D.

92

93   In order to consider a range of plausible ascertainment scenarios, we employ the general family-

94   based $k$-model of ascertainment (7). In its simplest form, this model stipulates that the probability

95   that a family is ascertained is proportional to $r^k$, where $k$ controls the model. For example, when

96   $k = 1$, the probability of ascertainment is strictly proportional to $r$: this is equivalent to classical

97   "single ascertainment". Similarly, when $k = 0$, so that every family with $r \geq 1$ is ascertained, this

98   model is equivalent to classical "complete" or "truncate" ascertainment.  We generalize this

99   model in two ways. First, we assume that ascertainment requires $r \geq 1$, that is, every ascertained

100  family contains at least one QI, but we allow that there may be additional preferential

101  ascertainment of families based on $t$ alone, that is, that investigators may preferentially ascertain

102  families with more affected individuals without knowing (or prior to knowing) the VOI status of

103  those additional individuals. Second, we allow that even an individual carrying the VOI may

104    develop disease due to any other independent causes at work in the general population. With

105    these two extensions in mind, our ascertainment model becomes

106         P[sibship is ascertained $\mid r, t] = c(r^k + t); for\ r\ \geq 1$, and 0 otherwise

107    where $c$ is a normalizing constant.

108

109    Let $f$ be the attributable penetrance, or the penetrance due to the VOI for HET individuals. (Note

110    that when $\gamma > 0$, $\beta$=P[AFF|HET]= $\gamma + f - \gamma f$. However, we focus here on estimation of $f$ itself

111    rather than $\beta$.) In what follows, we estimate $f$ in two ways:

112         (ii)    $\tilde{f}$ is obtained by counting the proportion of AFF individuals among all HET

113                individuals in the data set, after dropping one QI individual per family, that is,

114                applying the correction for single ascertainment;

115         (ii)    $\tilde{f}^*$ is obtained by counting the proportion of AFF individuals among all HET

116                individuals in the data set, that is, without applying any ascertainment correction.

117

118    $\tilde{f}^*$ is a naïve estimate, which would be correct if the families were not ascertained based on

119    either phenotype or genotype. It is, however, clearly incorrect under any of our ascertainment

120    models. Our interest in this estimate is to establish how biased it becomes under various

121    ascertainment scenarios. $\tilde{f}$ by contrast, does apply the frequently employed single ascertainment

122    correction, and again, our interest in $\tilde{f}$ is to establish how biased it will be under ascertainment

123    scenarios other than single ascertainment. Expected values of $\tilde{f}$ and $\tilde{f}^*$ were obtained via

124    simulation, by averaging each estimate's value across 1,000 replicates per generating condition,

125    and standard errors were obtained by averaging the standard deviation of each estimate across

126    those same 1,000 replicates. (While the expected values are easily calculated analytically, the

6

127   standard errors are not.) All simulations and calculations were done in MATLAB

128   (2021.9.10.0.1739362 (R2021a), Natick, Massachusetts: The MathWorks Inc.).

129

130   Let $s$ be the number of siblings in a family, and let N be the number of $s$-sized sibships in a

131   dataset. Fig 1 shows results for true single ascertainment (k=1), for $s = 2$, as a function of sample

132   size N. Here we assume that the true value of $f$=0.5. As can be seen, in this case, the mean of $\tilde{f}$ =

133   0.5, the generating value, as expected. But using $\tilde{f}^*$ the estimates are seriously upwardly biased

134   in all data sets, regardless of N. Note that because each sibship contains at least one QI, by

135   stipulation, the minimum value of $\tilde{f}^*$ is 0.50.

136

137   Note too that even the correct estimate $\tilde{f}$ shows considerable sampling variability. For instance,

138   with N=10, $\tilde{f}$ will be >70% or <30% in approximately 40% of all data sets when $f$ =50%. This

139   variability remains appreciable even for N=50.

140

141   For ascertainment models other than single, overall variability remains similar to what is shown

142   in Fig 1, but even $\tilde{f}$ tends to be biased, with mean $\tilde{f}$ = 0.60, 0.50, 0.43 and 0.38 for $k$ = 2, 1, 0

143   and −1, respectively. In all cases, the uncorrected $\tilde{f}^*$ will return even more biased estimates, with

144   mean $\tilde{f}^*$ = 0.89, 0.88, 0.87 and 0.86, for $k$ = 2, 1, 0 and −1, respectively.

145

146   Fig 2 shows the impact of the population prevalence $\gamma$ on average penetrance estimates.

147   Focusing first on single ascertainment ($k$=1) and $\tilde{f}$ =0.5, we can see that regardless of $k$, the

148   expected value of $\tilde{f}$ is relatively independent of $\gamma$ until $\gamma$ becomes quite high. Note that for $f$ =

149   0.5 and $\gamma = 0.5$, the actual probability that a VOI carrier is affected under our generating model

7

150     is $0.5 + 0.5 - (0.5)(0.5) = 0.75$, which is in line with the estimates returned by $\tilde{f}$. $\tilde{f}^*$ might be

151     said to be even more robust to $\gamma$, although this is because in this case $\tilde{f}^*$ is already close to the

152     top of the scale for $\gamma=0$. Moreover, $\tilde{f}^*$ appears not only robust to $\gamma$, but also to $f$ itself, with

153     estimates >70% even for $f=0.05$, and >80% for $f=0.05$ when $\gamma=0.5$.  These patterns repeat for

154     different values of $k$, with visible impact only on the magnitude of the bias for any given ($f, \gamma$)

155     combination. Ascertainment effects will be reduced as $s$ increases. Users who are interested in

156     investigating penetrance estimates for other ascertainment models, other combinations of

157     parameter values or other sibship sizes are encouraged to download the PenEst app:

158     https://github.com/MathematicalMedicine/PenetranceEstimator.

159

160     In general, our simulations show that under unsystematic ascertainment schemes, or in cases

161     where appropriate ascertainment corrections are not included in the estimation procedure, there

162     is a high risk of over-estimating the penetrance of any given VOI. This finding is consonant with,

163     and may in large part explain, reports for specific variants. For example, multiple coding variants

164     in *PRNP* had been reported to cause rare dominant monogenic neurodegenerative disease, but

165     there was a 30-fold higher prevalence of variants previously suggested to be causal in this gene

166     in ExAC compared to the expected frequency calculated from the estimated prevalence of the

167     disorder (8). Specifically for three variants the lifetime risk of developing disease was <10%.

168     Similarly, GWAS array data from the UK Biobank were used to estimate pathogenicity,

169     penetrance, and expressivity of putative disease-causing rare variants (MAF<1%) that were

170     directly genotyped and had good quality (9). Focused on maturity-onset diabetes of the young

171     and developmental disorders, many specific variants were found for which the penetrance --

172     estimated either in families ascertained for the presence of the VOI or in disease cohorts -- was

173    much higher than that obtained from a population-based cohort. These observations have

174    implications for genetic counselling, including the recommendation of invasive screening

175    procedures and administration of preventative treatment.

176

177    Some approaches to the interpretation of rare coding variants assume either full or high

178    penetrance (10), for the sake of simplicity. Extensive criteria have been proposed to claim a

179    causal relationship between variants and disease, and authors have urged caution in presuming

180    full penetrance for pathogenic variants (11). But in practice, penetrance remains an important

181    factor in assessing pathogenicity. For instance, the ACMGG/AMP joint consensus

182    recommendations (1) warns against ignoring the possibility of reduced penetrance in establishing

183    segregation of a VOI with a phenotype, but also instructs that "lack of segregation…provides

184    strong evidence against pathogenicity." (p. 15) And in practice, many laboratories will rule out

185    candidate VOIs when they are found among unaffected relatives. Particularly in the absence of a

186    rigorous and accurate estimate of the actual penetrance, this complicates the use of segregation

187    information in assessments of pathogenicity.

188

189    We close by noting that there is one essentially "ascertainment assumption free" (12) method for

190    estimating the penetrance, viz., by conditioning on all of the phenotypic data. This is the

191    ascertainment correction implicit in the usual LOD score (13-15), and also the LOD score

192    allowing for linkage disequilibrium or LD-LOD (6, 16, 17), and in principle any program that

193    allows calculation of the LOD score will support this method. As in Thompson (6) the

194    calculation is done here assigning the VOI (which plays the role of the "marker") and the disease

195    allele the same (rare) frequency (we have used 0.001 in the simulations), assuming complete

196    linkage disequilibrium between the two (D′ = 1), and also assuming 0 recombination between the

197    marker and the disease allele. Free parameters in the model are then the three penetrances; in our

198    calculations we also include the admixture parameter $\alpha$ of Smith (18), representing the

199    probability that any given family is of the "linked" type, which adds robustness when phenocopy

200    levels are high. Maximizing the LD-LOD over the free parameters gives us the LD-MOD, which

201    occurs at the maximum likelihood estimate (m.l.e.) of $\hat{f}$ of $f$ (12-15).

202

203    Fig. 3 shows results corresponding to the simulations in Fig 1A and Fig 2A, C. As can be seen, $\hat{f}$

204    behaves very much like $\tilde{f}$ when k = 1 (Fig 3A), but it retains almost complete robustness to

205    ascertainment, and also to $\gamma$ at least until $\gamma$ is quite large (Fig 3B). (As with $\tilde{f}$, as $\gamma$ gets very

206    large, $\hat{f}$ covers both cases due to the VOI and also cases among variant carriers due to other

207    causes.) Comparing Fig 3A with Fig 1A, $\hat{f}$ shows slightly greater sampling variability than $\tilde{f}$;

208    this is due to the inherent ascertainment correction built in to $\hat{f}$. The slight but systematic over-

209    or under-estimation of $f$ seen in Fig 3B is due to the small sample size; as N increases $\hat{f} \rightarrow$

210    $f$ (results not shown). However, in small samples the upward bias can be appreciable particularly

211    when $f$ is small; e.g., when $f = 0.05$ ($\gamma = 0$), for N = 20, the expected value of $\hat{f} = 0.165$.

212

213    Note, however, that while maximizing the LD-LOD is a highly ascertainment-robust method for

214    estimating $f$, the LD-MOD itself is not a good statistic for representing the strength of evidence

215    for co-segregation, because it is not additionally conditioned on ascertainment through the VOI.

216    However, once we ascertain so as to require the VOI to be present in the family, there is no

217    remaining LD information in the sibship, since LD information is conveyed entirely by the

218    marker allele frequencies in the parents. Therefore, we recommend using the ordinary (linkage

219    equilibrium) LOD, or LE-LOD, for assessing strength of evidence for co-segregation. Because

220    maximizing the LE-LOD itself will not return true m.l.e.s of $f$ under the LD model, we

221    recommend evaluating the LE-LOD at the maximizing model obtained from the LD-MOD, for a

222    statistic we annotate as LE-LOD(max). (This maximization procedure is not inherently

223    inflationary; see Supplemental Results (A). Thompson et al. (6) proposed a form of Bayes factor

224    for assessing evidence for co-segregation of the VOI with disease; see the Supplemental Results

225    (B) for some comparisons between their Bayes factor and LE-LOD(max).) Fig 3(D) shows the

226    distribution of the LE-MOD(max), for the same data shown in Fig 1. (Here parents are treated as

227    genotypically known but phenotypically unknown.) As expected, based on just a few 2-child

228    sibships, evidence of co-segregation of the VOI with disease is quite weak. It requires at least N

229    = 30 2-child families before there is a reasonable chance of obtaining a substantial LE-

230    LOD(max).

231

232    **Declaration of interests**

233    The authors declare no competing interests

234    **Acknowledgements:**

235    This work was supported in part by Mathematical Medicine LLC, with special thanks to Jo

236    Valentine-Cooper for creation of the PenEst app.

237

238    **Web resources**

239    PenEst app : https://github.com/MathematicalMedicine/PenetranceEstimator/

240

241    **Data and code availability**

242    https://github.com/MathematicalMedicine/PenetranceEstimator/

## References

1. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. Genet Med. 2015;17(5):405-24.

2. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. Nature. 2016;536(7616):285-91.

3. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. Nature. 2020;581(7809):434-43.

4. Backman JD, Li AH, Marcketta A, Sun D, Mbatchou J, Kessler MD, et al. Exome sequencing and analysis of 454,787 UK Biobank participants. Nature. 2021;599(7886):628-34.

5. Hodge SE, Vieland VJ. The essence of single ascertainment. Genetics. 1996;144(3):1215-23.

6. Thompson D, Easton DF, Goldgar DE. A full-likelihood method for the evaluation of causality of sequence variants from family data. Am J Hum Genet. 2003;73(3):652-5.

7. Ewens WJ, Shute NC. The limits of ascertainment. Ann Hum Genet. 1986;50(4):399-402.

8. Minikel EV, Vallabh SM, Lek M, Estrada K, Samocha KE, Sathirapongsasuti JF, et al. Quantifying prion disease penetrance using large population control cohorts. Sci Transl Med. 2016;8(322):322ra9.

9. Wright CF, West B, Tuke M, Jones SE, Patel K, Laver TW, et al. Assessing the Pathogenicity, Penetrance, and Expressivity of Putative Disease-Causing Variants in a Population Setting. Am J Hum Genet. 2019;104(2):275-86.

10. Jarvik GP, Browning BL. Consideration of Cosegregation in the Pathogenicity Classification of Genomic Variants. Am J Hum Genet. 2016;98(6):1077-81.

11. MacArthur DG, Manolio TA, Dimmock DP, Rehm HL, Shendure J, Abecasis GR, et al. Guidelines for investigating causality of sequence variants in human disease. Nature. 2014;508(7497):469-76.

269  12.     Ewens WJ, Shute NC. A resolution of the ascertainment sampling problem. I. Theory. Theor

270  Popul Biol. 1986;30(3):388-412.

271  13.     Greenberg DA. Inferring mode of inheritance by comparison of lod scores. Am J Med Genet.

272  1989;34(4):480-6.

273  14.     Elston RC. Man bites dog? The validity of maximizing lod scores to determine mode of

274  inheritance. Am J Med Genet. 1989;34(4):487-8.

275  15.     Vieland VJ, Hodge SE. The problem of ascertainment for linkage analysis. Am J Hum Genet.

276  1996;58(5):1072-84.

277  16.     Slager SL, Huang J, Vieland VJ. Power comparisons between the TDT and two likelihood-based

278  methods. Genet Epidemiol. 2001;20(2):192-209.

279  17.     Petersen GM, Parmigiani G, Thomas D. Missense mutations in disease genes: a Bayesian

280  approach to evaluate causality. Am J Hum Genet. 1998;62(6):1516-24.

281  18.     Smith CA. Testing for heterogeneity of recombination fraction values in human genetics. Ann

282  Hum Genet. 1963;27:175-82.

283  19.     Vieland VJ, Huang Y, Seok SC, Burian J, Catalyurek U, O'Connell J, et al. KELVIN: a software

284  package for rigorous measurement of statistical evidence in human genetics. Hum Hered. 2011;72(4):276-

285  88.

286     **Figure 1.** Swarm plots showing sampling distributions of penetrance estimates as a function of

287     number of families N.

288     Distributions of (A) $\tilde{f}$ and (B) $\tilde{f}^*$ are shown for simulations of 1000 replicates, with true

289     penetrance $f$=0.5. The number of sibs per family, s=2; phenocopy rate, $\gamma$=0. Users interested in

290     varying the parameters can use the PenEst app.

291

292     **Figure 2.** Expected values of penetrance estimates as a function of population prevalence $\gamma$ and

293     ascertainment parameter $k$.

294     Top row: expected values of (A) $\tilde{f}$ and (B) $\tilde{f}^*$ when the true penetrance $f$=0.5. Bottom row:

295     expected values of (C) $\tilde{f}$ and (D) $\tilde{f}^*$ when $f$=0.2 (lower line sets) or $f$=0.8 (upper line sets). The

296     number of sibs per family, s=2. Users interested in varying the parameters can use the PenEst

297     app.

298

299     **Figure 3.** (A) Swarm plots showing sampling distributions of $\hat{f}$, as obtained from maximizing

300     the LD-LOD, as a function of number of families N; (B) Expected values of $\hat{f}$ as a function of

301     population prevalence $\gamma$ and ascertainment parameter $k$, for $f = 0.2, 0.5$ and $0.8$, reading from

302     bottom to top of the plot, respectively; (C) Swarm plots showing the distribution of LE-

303     LOD(max) as a function of N. Data are the same as used to generate Figures 1 and 2,

304     respectively. All calculations were done using KELVIN (19).

Figure with three panels A, B, and C.

**A:** Plot of $\hat{f}$ versus $N$, with $N$ values at 1, 2, 4, 6, 8, 10, 30, 50. Horizontal dashed line at $\hat{f} = 0.5$.

**B:** Plot of $\hat{f}$ versus $\gamma$. Legend:
- $k=2$ (blue)
- $k=1$ (red)
- $k=0$ (yellow)
- $k=-1$ (purple)

Horizontal dashed lines at $\hat{f} = 0.8$, $0.5$, and $0.2$. The $\gamma$ axis shows values 0, 0.05, 0.1, and 0.5.

**C:** Plot of LE-LOD (max) versus $N$, with $N$ values at 1, 2, 4, 6, 8, 10, 30, 50.