Calculating variant penetrance using family history of disease and population data

<u>Authorship</u>

Thomas P Spargo¹, Sarah Opie-Martin¹, Cathryn M Lewis^{2,3}, Alfredo Iacoangeli^{1,4,5,*,#} and Ammar Al-Chalabi^{1,6,*,#}.

¹Maurice Wohl Clinical Neuroscience Institute, King's College London, Department of Basic and Clinical Neuroscience, London, UK;

²Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology & Neuroscience, King's College London, de Crespigny Park, London SE5 8AF, UK;

³Department of Medical and Molecular Genetics, Faculty of Life Sciences and Medicine, King's College London, London, UK;

⁴Department of Biostatistics and Health Informatics, King's College London, London, UK; ⁵NIHR Maudsley Biomedical Research Centre (BRC) at South London and Maudsley NHS Foundation Trust and King's College London, London, UK;

⁶King's College Hospital, Bessemer Road, London, SE5 9RS, UK.

*co-senior author

[#]correspondence should be addressed to <u>alfredo.iacoangeli@kcl.ac.uk</u> and <u>ammar.al-</u> <u>chalabi@kcl.ac.uk</u>

Abstract

Genetic penetrance is the probability of a phenotype manifesting given that one harbours a specific variant. For most Mendelian genes, penetrance is high, but not complete, and may be age-dependent. Accurate estimates of penetrance are important in many biomedical fields including genetic counselling, disease research, and for gene therapy. The main methods for its estimation are limited in situations where large family pedigrees are not available, the disease is rare, late onset, or complex. With the advance of high-throughput technologies, population-scale genetic data is available for an increasing range of genetic diseases. Here we present a novel method for penetrance estimation in autosomal dominant phenotypes. It uses population-scale data regarding the distribution of a variant among unrelated people affected and unaffected by an associated phenotype and can be restricted to samples of affected people by considering family disease history. The approach avoids kinship-specific penetrance estimates and the ascertainment biases that can arise when sampling rare variants among control populations. We test the method upon candidate variants and diseases, demonstrating that our estimates align with those derived using established methods. We have implemented the method in a public web server (https://adpenetrance.rosalind.kcl.ac.uk) and made it available as an open-source R library (https://github.com/ThomasPSpargo/adpenetrance).

Key words: Penetrance, mutation, autosomal dominant, disease, familial, sporadic

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

Introduction

Penetrance is the probability of developing a specific trait given that a person harbours a certain genetic variant or set of variants. Some pathogenic variants are fully penetrant, and people harbouring them always develop the associated phenotype. For instance, a trinucleotide CAG repeat expansion within the *HTT* gene is fully penetrant for Huntington's Disease by 80 years of age among people harbouring an expansion variant larger than 41 repeated CAG units (1). For many variants however, penetrance is incomplete, and those with risk variants can remain unaffected throughout their life. For example, the p.Gly2019Ser variant of the *LRRK2* gene exhibits incomplete penetrance for Parkinson's Disease (PD), meaning that it elevates risk for PD but does not necessarily result in its manifestation (2).

In medical genetics, estimating the penetrance of a given variant or set of variants is important for the correct interpretation of genetic test results, something that will be increasingly valuable as genome sequencing becomes routine, both within and outside clinical practice. With the advance of precision medicine and gene therapy, being able to accurately estimate the penetrance of a large spectrum of human genetic variants is crucial (3-6).

There are several existing methods for penetrance estimation. The first and most widely used is based on the statistical examination of how the variant segregates with the phenotype within pedigrees (7). However, the generalisability of estimates derived from specific families may be limited. Other approaches involve examination of the incidence of disease in a sample of unrelated people who harbour a variant (8, 9). Without systematic sampling, these estimates can be affected by ascertainment bias. Where large pedigrees are not available, or if disease is rare or late onset, these techniques may not be possible (10).

Estimating penetrance for a variant of unknown significance identified, for example, as a result of genome sequencing-based screening can be particularly challenging. The problem is exemplified by the large number of reported *SOD1* gene variants in amyotrophic lateral sclerosis (ALS): although *SOD1* variants are cumulatively one of the most common causes of ALS, over 180 ALS-associated variants in the gene are reported to date (11, 12). Family pedigrees suitable for establishing penetrance are available for only a minority of these.

We have developed a new method to calculate penetrance for variants with an autosomal dominant inheritance pattern using population level data from unrelated people who are and are not affected by the associated phenotype (case and control populations). It can be operated using variant information drawn only from affected populations, stratified according to family history between 'familial' and 'sporadic' disease presentations. This approach is based on our previously published model of disease which explains how variant penetrance and sibship size determine the presentation or absence of a disease for families in which the variant occurs (13).

The method is complementary to, and fills an important gap left by, existing techniques. Using population-scale data, it takes full advantage of the rapidly growing quantity of genetic data that are being generated for a wide range of human disease and, therefore, it is ideally placed to be a valuable tool in the precision medicine era. Moreover, the capacity to assess penetrance based on the distribution of a variant between samples of unrelated people drawn only from the affected population allows estimates unbiased by kinship-specific effects or ascertainment of unaffected population members.

We have tested the approach in four variant-disease case examples, drawing upon the most common and widely studied autosomal dominant risk variants for each disease: the p.Gly2019Ser variant of the *LRRK2* (OMIM: 609007) gene for PD (2); variants in the *BMRP2* gene (OMIM: 600799) for heritable pulmonary arterial hypertension (PAH) (14); and variants in the *SOD1* (OMIM: 147450) and *C9orf72* (OMIM: 614260) genes, for ALS (15, 16).

Methods

Model

The disease model our method builds upon (13) makes the following assumptions: a rare dominant variant is necessary but not sufficient for disease to occur, therefore penetrance is not complete and people within a family who do not harbour the variant are not affected; all individuals harbouring the variant are ascertained; all variants are inherited from exactly one parent, thus there are no homozygous carriers or *de novo* variants.

The model calculates three probabilities for a nuclear family where one parent harbours a given variant: that no family members are affected, P(unaffected); that exactly one member is affected, P(sporadic); and that more than one member is affected, P(familial). These probabilities are determined by penetrance, f, and sibship size, N. In a family with N siblings:

$$P(unaffected) = (1-f)\left(1-\frac{f}{2}\right)^{N}$$
(1)

where the parent is unaffected, and none of the sibs are affected (each being transmitted the high-risk variant with probability $\frac{1}{2}$).

$$P(sporadic) = f\left(1 - \frac{f}{2}\right)^{N} + N\left(\frac{f}{2}\right)\left(1 - \frac{f}{2}\right)^{N-1}(1 - f)$$
(2)

where either the parent is affected and no siblings are affected, or the parent is unaffected and exactly one of the sibs is affected. Then,

$$P(familial) = 1 - P(unaffected) - P(sporadic)$$
$$= 1 - \left(\left(1 - \frac{f}{2}\right)^{N} + N\left(\frac{f}{2}\right) \left(1 - \frac{f}{2}\right)^{N-1} (1 - f) \right)$$
(3)

Application to penetrance calculation

Conversely, given the observed rates of the (unaffected, sporadic, familial) disease states in families where the variant occurs and the average sibship size for these families, we can estimate penetrance. We can also estimate penetrance based on the observed rates of families presenting as unaffected versus 'affected', a fourth disease state whereby P(affected) = P(familial) + P(sporadic). Observed disease state rates can be derived as a weighted proportion of estimates of heterozygous variant frequency given for people across a valid subset of the four defined states (see table 1); the appropriate weighting factors will vary based on the disease states for which variant frequencies are given. Sibship size can be estimated for the sample either directly, based on the average sibship size among the described families, or indirectly, by designating an estimate representative of the sampled population (e.g. available within global databases).

medRxiv preprint doi: https://doi.org/10.1101/2021.03.16.21253691; this version posted March 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

perpetuity. It is made available under a CC-BY 4.0 International license .

Variant frequencies provided	Required weighting factors
Familial (M_F) ,	$W_F = P(F A),$
Sporadic (M_S)	$W_S = P(S A)$
Familial (M_F) ,	$W_F = P(F A) \times P(A),$
Unaffected (M_U)	$W_U = 1 - P(A)$
Sporadic (M_S) ,	$W_S = P(S \mathbf{A}) \times P(\mathbf{A}),$
Unaffected (M_U)	$W_U = 1 - P(\mathbf{A})$
Familial (M_F) ,	$W_{\rm F} = P(F {\rm A}) \times P(A),$
Sporadic (M_S),	$W_S = P(S A) \times P(A),$
Unaffected (M_U)	$W_U = 1 - P(A)$
Affected (M_A) ,	$W_A = P(A),$
Unaffected (M_U)	$W_{U} = 1 - P(A)$

Table 1. Valid disease state combinations and weighting factors used to estimate disease state rates associated with a given variant as described in Figure 1 and the supplementary methods. $M_{F,S,U,A}$ = variant frequencies in the familial, sporadic, unaffected, and affected states; $W_{F,S,U,A}$ = weighting factors for the familial, sporadic, unaffected, and affected states; P(A) = the probability of a member of the sampled population being affected; P(F|A) = disease familiality rate; P(S|A) = disease sporadic rate.

Our method involves three operations and an optional further step for deriving error in the estimate. These processes are summarised as a flowchart in Figure 1 and outlined in detail in the supplementary methods. In this approach, we assume that: in variant frequency estimates, disease state classifications are assigned according to the status of the sampled person and first-degree relatives only; individual families are represented only once in variant frequency estimates; weighting factors and average sibship size represent absolute values; the value specified for sibship size is representative of sibship size across disease state groups; in families where the variant occurs, the associated trait can only manifest owing to that variant.



medRxiv preprint doi: https://doi.org/10.1101/2021.03.16.21253691; this version posted March 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in

perpetuity. It is made available under a CC-BY 4.0 International license

Figure 1. Flowchart summarising the key operations within this penetrance estimation approach. Operation 1: variant frequencies (M) and weighting factors (W) are defined for a valid subset of the familial (F), sporadic (S), unaffected (U), and affected (A) states to calculate rate of one of these states, arbitrarily labelled state X, among families in which the variant occurs drawn from those states for which data were provided. This is the observed rate of state X, $R(X)^{obs}$, and Table 1 summarises all valid state combinations. Operation 2: Equations 1-3 are applied to calculate P(familial), P(sporadic), P(unaffected), and P(affected) at a specified sibship size (N) and for a series of penetrance values, $f_i = 0, ..., 1$. The rate of arbitrary state X that is expected at each f_i among variant harbouring families from those states represented in Operation 1, $R(X)_i^{ex}$, is calculated and stored alongside the corresponding f_i in a lookup table. Operation 3: The lookup table is queried using $R(X)^{obs}$ to identify the closest $R(X)_i^{ex}$ value and then the f_i to which this corresponds. This f_i is taken as the penetrance estimate. Optional step: upper and lower confidence intervals for $R(X)^{obs}$ can be calculated from error in the given estimates of M (17). Penetrance is estimated as in Operation 3 for the bounds of these intervals.

Tool access

We have made this method available as an R function (R version 3.6.1) and, leveraging the R Shiny package (version 1.4.0.2), also developed a publicly available web resource (https://adpenetrance.rosalind.kcl.ac.uk) that facilitates easy use of the method. The source code of the R library is available on GitHub

(https://github.com/ThomasPSpargo/adpenetrance). The web tool is further described in the supplementary methods and Figure 2 presents an example of its usage.

	Penetrance calculato	<u>or</u>		
Disease states represented in data:	Data format:	Include error propagation?		
	 Variant counts with sample size 	Provide standard errors		
Affected	states			
Please provide data for any combination of two or affected' state represents families in whom at leas Familial parameters: Variant frequency estimate	more of the 'familial', 'sporadic' and 'unaffected t one person has developed disease; it is the su Sporadic parameters: Variant frequency estimate	d' disease states OR the 'affected' and 'unaffected' states. The familial and sporadic disease states.		
0.148	0.012			
Lower confidence interval	Lower confidence interval			
0.115	0.007			
Confidence level	Confidence level			
	05%	-		
Define weighting factors:	Sibship size:	Sibship data repository		
Cases: Familial disease rate	1.54344477628159	Total Fartility Pata (Marid bank 2020)		
0.05		Total Fertility Rate (world bank, 2020):		
Cases: Sporadio disease rate		European Union		
0.95	Confidence level for penetrance estimate:	Select year		
	95%	2018 •		
		Set sibship size		
Table 1. Observed disease state rate and correspo	Calculate Disease modelled	state probabilities expected across values of penetran for a population of sibship size 1.543		
Lower CI Estimate	Upper Standard CI error			
Observed familial 0.281 0.394 rate	0.506 0.058 0.8	Observed ra		
Expected familial 0.281 0.394 rate	0.506 NA	— Familial Sporadic		
Penetrance 0.494 0.660	0.812 NA to 0.4	State probab		
Table 2. Disease state probabilities at the penetran estimate	ce	Unaffected		
Table 2. Disease state probabilities at the penetran estimate Lower CI Estimate Upper	CI 0.2	Unaffected		
Table 2. Disease state probabilities at the penetranestimate Lower CI Estimate Familial 0.189 0.321 0.46	CI 0.2	Unaffected		

Figure 2, example interface and output of the ADPenetrance web tool

(https://adpenetrance.rosalind.kcl.ac.uk). Here we show the example of penetrance of SOD1 variants for amyotrophic lateral sclerosis in a European population, applying variant frequency estimates for familial and sporadic ALS patients of European ancestry and the average Total Fertility Rate for the European Union in 2018 (18, 19).

Case examples

Input parameters for the presented case studies have been estimated using publicly available data. We estimated variant frequency in the familial, M_F , and sporadic, M_S , states in all cases and, in case 1, additionally ascertained this for the unaffected state, M_U , from control samples. In all cases, we derived the standard error of these values, $\sigma_{\overline{M_X}}$, to allow for assessment of error in the penetrance estimate. Variant frequency estimates were weighted to calculate the observed rate of arbitrary state X, $R(X)^{obs}$, among variant harbouring families using the factors presented in Table 1. Accordingly, the frequencies of familial, P(F|A), and sporadic, P(S|A), disease among the affected population, A, were defined in all cases, based on the first-degree familiality rates of that trait; note that P(S|A) = 1 - P(F|A). The probability of a population member being affected, P(A), was defined for case 1 only, using estimates of lifetime risk for that trait.

Sibship size, N, was estimated in each case based on the Total Fertility Rates reported in the World Bank database (19) for the world region(s) which best represent the population from which variant frequency estimates were drawn.

An R script detailing the calculations made for each case study can be found within our GitHub repository.

Case 1: LRRK2 penetrance for PD

We estimated the penetrance of the p.Gly2019Ser variant of *the LRRK2* gene for PD. This case was used to illustrate the flexibility of this method for application using data drawn from several combinations of the defined disease states.

The first-degree familiality rate of PD, about 0.105, was used to estimate P(F|A) and P(S|A) (Ref. 20, 21). P(A) was estimated as 1 in 37 (0.027), the lifetime risk of developing PD (22).

We estimated $M_{F.S.U}$ based on a report which aggregates data from 24 world populations (23). Of 5,123 unrelated people with familial PD manifestations, 201 ($M_F = 0.039$, $\sigma_{\overline{M_X}} = 2.71 \times 10^{-3}$) harboured the *LRRK2* p.Gly2019Ser variant, compared to 179 of 14,253 with sporadic PD manifestations ($M_S = 0.013$, $\sigma_{\overline{M_S}} = 9.33 \times 10^{-4}$) and 11 of 14,886 unaffected controls ($M_U = 7.39 \times 10^{-4}$, $\sigma_{\overline{M_U}} = 2.23 \times 10^{-4}$).

This intercontinental cohort largely describes people from European, North American and Asian countries but there is no single predominant region. We accordingly estimated that N = 1.646 by aggregating Total Fertility Rate estimates available in the World Bank database (19) across each of the reported population samples. For each population, this was weighted by its proportional contribution to the total sample; see Table S1 for further details.

Case 2: BMPR2 penetrance for heritable PAH

We estimated the penetrance of variants in the *BMPR2* gene for heritable PAH, a gene for which the low penetrance of pathogenic variants is well established (24).

Input parameters were defined based on only people with idiopathic (sporadic) or heritable PAH diagnoses (14). This captures people with and without family disease history and

excludes those with PAH manifestations associated with comorbidities or drug exposure.

We estimated P(F|A) and P(S|A) using the first-degree familiality rate of heritable PAH, about . 055 of people affected by either idiopathic or familial PAH (24).

In this case, to minimise any study specific bias, we applied data from two reports to build independent estimates for each of $M_{F,S}$.

The first dataset (14), presents a moderately large sample of people with familial and sporadic PAH. Of 247 people with familial PAH, 202 harboured *BMPR2* variants ($M_F = 0.818, \sigma_{\overline{M_F}} = 0.025$), compared to 200 of 1174 in the sporadic state ($M_S = 0.170, \sigma_{\overline{M_S}} = 0.011$). It is possible that this data may violate two assumptions of our approach. First, information on familial clustering was reportedly unavailable and so some families may be represented more than once in the familial state. Second, it is not specified whether disease familiality is defined only by the disease status of first-degree relatives.

The second dataset (25), overcomes a limitation of the first as each family is represented only once in variant counts. However, the sample is smaller in size. It is reported that 40 of 58 people with familial PAH ($M_F = 0.690, \sigma_{\overline{M_F}} = 0.061$) harboured *BMRP2* variants, compared to 26 of 126 in the sporadic state ($M_S = 0.206, \sigma_{\overline{M_S}} = 0.036$). Variant counts are additionally reported separately for small genetic variations (point mutations and indels) and large genetic rearrangements in *BMPR2*, which allowed penetrance estimation stratified by variant type. It is not reported whether disease states are defined according to the status of first-degree relatives only.

The first cohort samples people from Asian, European, and North American populations; French, German and Italian cohorts comprise about 60% of the sample (14). The second cohort samples people exclusively from Western Europe (25). We therefore estimated that N = 1.543 in both instances, the Total Fertility Rate of the European Union in 2018 (Ref. 19).

Cases 3 and 4: SOD1 and C9orf72 penetrance for ALS

We estimated the penetrance of variants in the *SOD1* and *C9orf72* genes for ALS. In *SOD1*, we examined the aggregated of penetrance of various *SOD1* variants harboured by people with ALS. For *C9orf72*, we examined the penetrance of a single pathogenic variant, a hexanucleotide GGGGCC repeat expansion. These penetrances have been historically difficult to establish without incurring kinship-specific biases. It is an ideal candidate for usage of our method.

The first-degree familiality rate of ALS, about 0.050, was applied to define P(F|A) and P(S|A) in these cases (26, 27).

We drew upon the results of two recent meta-analyses to estimate $M_{F,S}$ for SOD1 and C9orf72 (18, 28). As variant frequencies differed between Asian and European ancestries, we model these, and therefore penetrance, separately for each group. We derive $\sigma_{\overline{M_{F,S}}}$

using z-score conversion from the 95% confidence intervals (95% CIs) reported: for the arbitrary state X,

$$\sigma_{\overline{M_X}} = \frac{M_X - M_X^{95\% lower}}{Z} \tag{4}$$

where z = 1.96 and $M_X^{95\% lower}$ is the lower 95% CI bound of the estimate M_X .

Accordingly, we identified that, in Asian ALS populations: SOD1 variants were harboured by 0.300 ($\sigma_{\overline{M_F}}$ = 0.025) of people with familial and 0.015 ($\sigma_{\overline{M_S}}$ = 2.55 × 10⁻³) with sporadic disease; the C9orf72 repeat expansion was harboured by 0.04 ($\sigma_{\overline{M_F}} = 0.010$) of people with familial and 0.01 ($\sigma_{\overline{M_S}} = 5.10 \times 10^{-3}$) with sporadic disease. In European ALS populations: SOD1 variants were harboured by 0.148 ($\sigma_{\overline{M_F}} = 0.017$) of people with familial and 0.012 ($\sigma_{\overline{M_s}} = 2.55 \times 10^{-3}$) with sporadic disease; the *C9orf72* repeat expansion was harboured by 0.32 ($\sigma_{\overline{M_F}}$ = 0.020) of people with familial and 0.05 ($\sigma_{\overline{M_S}}$ = 5.10 × 10⁻³) with sporadic disease.

In these datasets, the Asian ancestry cohorts were predominantly individuals from East Asia, with small proportion from South Asia. The European ancestry cohorts primarily comprise people from European countries, with some from North America and Australasia. Accordingly, N was estimated for the Asian population samples as 1.823, the Total Fertility Rate for East Asia and Pacific in 2018, and for the European population as 1.543, the Total Fertility Rate for the European Union in 2018 (19).

Results

Here we summarise the input data and results of the case studies modelled (see table 2).

In case 1, we estimated the penetrance of the p.G2019S variant of the *LRRK2* gene for PD, taking estimates of $M_{\rm F,S,U}$ to allow estimation via four of the five possible disease state combinations presented in table 1. The output estimates (see table 2) were consistent across the modelled disease state combinations, with some discordance between estimates derived with and without the inclusion of the unaffected disease state. We expect this to reflect that the variant is rare in the unaffected (control) population and the variant frequency estimate may be affected by an ascertainment bias.

In case 2, we estimate the penetrance of variants in *BMPR2* for PAH, drawing estimates of $M_{\rm F,S}$ from two distinct reports (see table 2). For the first sample (14), we found penetrance of 0.395 (95% *CI*: 0.356, 0.433), compared to 0.303 (95% *CI*: 0.211, 0.390) for the second sample set (25), in which penetrance was comparable between the defined *BMPR2* variant subtypes. The marginally higher penetrance estimate observed for first dataset reflects differences observed in $M_{F,S}$ between the cohorts and may be affected by unspecified family clustering within this sample set. It is not known for either dataset whether family history classifications were restricted to first-degree relatives only and so the estimates obtained may be slightly inflated. With the available data these possibilities cannot be explored further.

In cases 3 and 4, we estimated the penetrance of variants in *SOD1* and *C9orf72* for ALS, drawing estimates of $M_{\rm F,S}$ for each gene in Asian and European populations separately based on the findings of recent meta-analyses (see table 2). We found the penetrance of *SOD1* variants to be 0.749 (95% *C1*: 0.629, 0.864) in Asian and 0.660 (95% *C1*: 0.494, 0.812) in European populations, and the penetrance of the pathogenic *C9orf72* hexanucleotide repeat expansion to be 0.282 (95% *C1*: 0.023, 0.514) in Asian and 0.449 (95% *C1*: 0.377, 0.518) in European populations. These estimates demonstrate consistency within genes across populations and indicate that the penetrance for ALS is greater in people harbouring *SOD1* variants than in those harbouring the *C9orf72* expansion. Table S2 presents additional penetrance estimates made for widely-described SOD1 variants: penetrance was estimated to be 0.917 for p.Ala5Val, 0.617 for p.Ile114Thr, and 0.0009 for p.Asp91Ala.

Case study	Data subset	Variant frequency in familial state (standard error)	Variant frequency in sporadic state (standard error)	Variant frequency in unaffected state (standard error)	Lifetime risk of disease [§]	Proportion familial [*]	Average sibship size [†]	States modelled #	Familial disease rate among those harbouring the variant across states modelled (95% confidence interval)	Penetrance (95% Confidence interval)
-	-	$M_F(\sigma_{\overline{M_F}})$	$M_{S}\left(\sigma_{\overline{M_{S}}} ight)$	$M_U(\sigma_{\overline{M_U}})$	P(A)	P(F A)	Ν	-	R(X)	f
<i>LRRK2</i> <i>p.G2019S</i> for PD (Ref. 23)	-	0.039 (2.71x10 ⁻³)	0.013 (9.32x10 ⁻⁴)	7.39x10 ⁻⁴ (2.23x10 ⁻⁴)	0.027	0.105	1.646 ^{†a}	F, S, U	0.098 (0.059, 0.137)	0.336 (0.258, 0.401)
								F, S	0.268 (0.229, 0.307)	0.453 (0.394, 0.511)
								F, U	0.134 (0.064, 0.205)	0.306 (0.221, 0.368)
								S, U	0.297 (0.170, 0.424) ^ø	0.196 (0.103, 0.306)
BMRP2 variants for PAH	All variants (Ref. 14)	0.818 (0.025)	0.170 (0.011)	-	-	0.055	1.543 ⁺ c	F, S	0.218 (0.195, 0.218)	0.395 (0.356, 0.433)
	All variants (Ref. 25)	0.690 (0.061)	0.206 (0.036)	-	-	0.055	1.543 ⁺	F, S	0.163 (0.111, 0.215)	0.303 (0.211, 0.390)
	Small variants (Ref. 25)	0.569 (0.065)	0.159 (0.033)	-	-	0.055	1.543 ⁺	F, S	0.173 (0.107, 0.238)	0.319 (0.204, 0.427)
	Large variants (Ref. 25)	0.121 (0.043)	0.048 (0.019)	-	-	0.055	1.543 ^{+c}	F, S	0.129 (0.012, 0.246)	0.243 (0.023, 0.439)
SOD1 variants	Asian	0.300 (0.025)	0.015 (2.55×10 ⁻³)	-	-	0.050	1.823 ^{†b}	F, S	0.513 (0.420, 0.606)	0.749 (0.629, 0.864)
for ALS (Ref. 18)	European	0.148 (0.017)	0.012 (2.55×10 ⁻³)	-	-	0.050	1.543 ^{+c}	F, S	0.394 (0.281, 0.506)	0.660 (0.494, 0.812)
C9orf72 ^{RE}	Asian	0.04 (0.010)	0.01 (5.10×10 ⁻³)	-	-	0.050	1.823 ^{+b}	F, S	0.174 (0.013, 0.335)	0.282 (0.023, 0.514)
for ALS (Ref. 28)	European	0.32 (0.020)	0.05 (5.10×10 ⁻³)	-	-	0.050	1.543 ⁺ c	F, S	0.252 (0.208, 0.296)	0.449 (0.377, 0.518)

Table 2. Penetrance estimation for the present case studies. [§]Lifetime disease risk is only required as a weighting factor where the unaffected (control) population are represented within the data given (see Table 1); ^{*}Proportion sporadic is defined as 1 - proportion familial (P(S|A) = 1 - P(F|A)); [†]Estimated using Total Fertility Rates reported for the: populations sampled to calculate variant frequencies (see Table S1) ^{†a}, East Asia and Pacific^{†b}, or European Union^{†c} regions in 2018 (Ref. 19); [#]F=familial, S=sporadic, U=unaffected (controls); ^øRate of sporadic disease has been calculated here because the familial state is not represented; C9orf72^{RE} = the pathogenic C9orf72 GGGGCC hexanucleotide repeat expansion.

Discussion

We have developed a novel approach to estimate the penetrance of genetic variants which confer risk for autosomal dominant traits. The method was tested via application to several variant-disease case studies.

Our penetrance estimates of the *LRRK2* p.G2019S variant for PD closely matched those previously obtained when analysing data that is not liable to inflation owing to selection of familial cases (2). Such studies make lifetime penetrance estimates for this variant between 0.24 (95% *CI*: 0.135, 0.437) and 0.45 (CI not reported).

Our estimates for the penetrance of *BMRP2* variants for PAH also aligned with existing estimates. Longitudinal analysis of disease trends among 53 families harbouring *BMRP2* variants finds penetrance as 0.27 overall, 0.42 for women and 0.14 for men (29). Our slightly higher estimate could reflect that a broader definition of familiality was used in the assessed samples, however this cannot be tested with the data available.

The estimates we generated in the *SOD1* and *C9orf72* case studies aligned with current understanding of the penetrance of variants in these genes for ALS.

For SOD1 variants, penetrance for ALS in a normal lifespan is reportedly incomplete and differs between individual variants (10, 30). The widely-described p.Ala5Val (formally p.Ala4Val) variant has been recorded to have penetrance of .91 by age 70 (31). Among other variants, penetrance is less apparent and can be expected to be lower than this (10, 30). Of the best characterised variants, p.Ile114Thr approaches complete penetrance in some, but not all, pedigrees and p.Asp91Ala reaches polymorphic frequency in some populations, with ALS typically arising with an autosomal recessive pattern (10, 11, 31). We drew estimates which align with these observations when modelling penetrance of the heterozygous forms of these three variants individually (see Table S2), estimating it to be 0.917 for p.Ala5Val, 0.617 for p.lle114Thr, and 0.0009 for p.Asp91Ala. These findings highlight the spectrum of penetrance across variants in SOD1. Our estimate for the p.Asp91Ala variant in particular is compatible with and supports the hypothesis that it is associated with ALS via a recessive or oligogenic inheritance pattern (32). The absence of p.Asp91Ala within the familial ALS database sampled further corroborates this finding. Accordingly, our penetrance estimates in Asian and European populations can be taken to suitably represent an aggregated penetrance of risk variants in SOD1 for ALS; some variation between populations can be expected, reflecting differences in the admix of variants between them.

For *C9orf72*, we modelled the penetrance of its pathogenic hexanucleotide repeat expansion for ALS. Pleiotropy is a well-established characteristic of this variant, additionally conferring risk for frontotemporal dementia and, to a lesser degree, other neuropsychiatric conditions (33). In people who harbour the variant, age-dependent penetrance for ALS and frontotemporal dementia is about equal and has been reported as almost complete at around age 80 (34). This estimate is however liable to inflation from biased ascertainment of affected people, and unaffected people are observed to harbour this variant more often than would be expected if it were accurate (16, 33, 34). Adjusted for possible ascertainment bias, the penetrance for either ALS or frontotemporal dementia is tentatively reported as 0.90 by age 83. Accounting for lifetime risk of each phenotype and their respective

familiality rates, people of European ancestry who harbour the *C9orf72* repeat expansion appear to develop ALS or frontotemporal dementia with comparable frequency, with 1.012 cases of ALS emerging per case of frontotemporal dementia (See Table S3; 28, 35-38). It is therefore reasonable to predict that penetrance of this variant for ALS would be around 0.45, a value comparable to our findings.

The method we present has high validity. Criterion and face validity are shown across the penetrance estimates outlined in the present paper, aligning with those made using other techniques and current understanding of the assessed cases. Construct validity is also demonstrated. In the ALS case studies, we found disease risk to be greater for those harbouring a pathogenic *SOD1* variant than for those with the *C9orf92* repeat expansion. This aligns with the multi-step model of ALS, where harbouring *SOD1* variants is associated with a 2-step disease process, converse to the 3-step process associated with the *C9orf72* repeat expansion (39).

The data necessary to operate the present approach is distinct from that of prior penetrance estimation techniques which examine patterns of disease among affected people, allowing it to be assessed in unrelated populations rather than families. The estimates are therefore unaffected by kinship-specific modifiers and are instead applicable to the region from which data are drawn.

Where the analysis is confined to people affected by disease, across the familial and sporadic states, we circumvent the ascertainment biases affecting designs which examine the distribution of a variant between affected and unaffected populations (9). In instances where analysis includes data for unaffected samples (i.e. controls harbouring the variant) these would not be avoided; ascertainment of controls compared to cases has equivalent challenges irrespective of the penetrance estimation approach. However, as our method does not require this information if data of familial and sporadic cases are available, it does not majorly limit the approach.

Furthermore, limitations of ascertainment will diminish as huge datasets of genetic and phenotypic information available within public databases become increasingly available. Therefore, the usefulness of penetrance estimates generated through population data will grow as the size and scope of genetic data held in such datasets expands, facilitating accurate estimation how of disease manifestations are distributed within the population in relation to harboured genetic variation (9).

A limitation of this approach is the definition of familiality, which we have defined as the occurrence of the studied trait in a first-degree relative. In practice, familial disease may be defined using various criteria, for example considering the disease status of second- or third-degree relatives, or including related diseases that may share a genetic basis (27, 40). For example, ovarian and breast cancer, or ALS and frontotemporal dementia each share a genetic basis, and it is reasonable to consider a family history of frontotemporal dementia when assessing familiality in a person with ALS. If the extended kinship is incorporated within familial disease state definitions, then the familial rate will trend upwards and inflate penetrance estimates. However, the use of a wider definition of being affected is acceptable, although it will yield penetrance estimates for the joint condition.

This method is suitable for calculating the point, rather than age-dependent, penetrance of a variant. It can be applied to derive penetrance for an individual variant or for an aggregated set of variants, with the latter indicating an averaged burden of variants meeting the given criteria. It can be applied to any form of germline genetic variation that is associated with a given trait via an autosomal dominant inheritance pattern.

In a scenario where penetrance can be estimated via multiple approaches, we recommend that researchers use each method available to them, given the complimentary nature of these techniques. If the results of multiple approaches conflict, we would suggest inspection of the suitability of the input data given for each method and to prioritise the result obtained from the method which this fits best.

In conclusion, our novel method for penetrance estimation fills an important gap in medical genetics because, making use of the available amounts of population-scale data, it enables the unbiased and valid calculation of penetrance in genetic disease instances that would be otherwise difficult or impossible using existing methods. It serves to expand the range of genetic diseases and variants for which high-quality penetrance estimates can be obtained, as we illustrate in the ALS case examples. Estimates drawn via this approach have clear clinical utility and will be useful for guiding the interpretation of genetic test results that reveal an individual to harbour a characterised risk variant. They have wider relevance to the population than those obtained by studying particular kinships and will be more interpretable for clinical professionals.

The tool code is available on GitHub (<u>https://github.com/ThomasPSpargo/adpenetrance</u>) and the method is available and free to use via a public webserver (<u>https://adpenetrance.rosalind.kcl.ac.uk</u>).

Conflict of interest statement

CML sits on the SAB for Myriad Neuroscience.

Acknowledgements

This is an EU Joint Programme-Neurodegenerative Disease Research (JPND) project. The project is supported through the following funding organizations under the aegis of JPNDhttp://www.neurodegenerationresearch.eu/ (United Kingdom, Medical Research Council MR/L501529/1 to A.A.-C., principal investigator [PI] and MR/R024804/1 to A.A.-C., PI]; Economic and Social Research Council ES/L008238/1 to A.A.-C. [co-PI]) and through the Motor Neurone Disease Association. This study represents independent research partly funded by the National Institute for Health Research (NIHR) Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's College London. The work leading up to this publication was funded by the European Community's Horizon 2020 Programme (H2020-PHC-2014-two-stage; grant 633413). We acknowledge use of the research computing facility at King's College London, Rosalind (https://rosalind.kcl.ac.uk), which is delivered in partnership with the National Institute for Health Research (NIHR) Biomedical Research Centres at South London & Maudsley and Guy's & St. Thomas' NHS Foundation Trusts and part-funded by capital equipment grants from the Maudsley Charity (award 980) and Guy's and St Thomas' Charity (TR130505). The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR, King's College London, or the Department of Health and Social Care.

References:

Langbehn DR, Brinkman RR, Falush D, Paulsen JS, Hayden MR. A new model for 1. prediction of the age of onset and penetrance for Huntington's disease based on CAG length. Clin Genet. 2004; 65(4):267-77. doi:10.1111/j.1399-0004.2004.00241.x

2. Goldwurm S, Tunesi S, Tesei S, Zini M, Sironi F, Primignani P, et al. Kin-cohort analysis of LRRK2-G2019S penetrance in Parkinson's disease. Mov Disord. 2011; 26(11):2144-5. doi:10.1002/mds.23807

Dewey FE, Grove ME, Pan C, Goldstein BA, Bernstein JA, Chaib H, et al. Clinical 3. Interpretation and Implications of Whole-Genome Sequencing. JAMA. 2014; 311(10):1035-45. doi:10.1001/jama.2014.1717

4. Kalia SS, Adelman K, Bale SJ, Chung WK, Eng C, Evans JP, et al. Recommendations for reporting of secondary findings in clinical exome and genome sequencing, 2016 update (ACMG SF v2.0): a policy statement of the American College of Medical Genetics and Genomics. Genet Med. 2017; 19(2):249-55. doi:10.1038/gim.2016.190

5. Saelaert M, Mertes H, Moerenhout T, De Baere E, Devisch I. Criteria for reporting incidental findings in clinical exome sequencing – a focus group study on professional practices and perspectives in Belgian genetic centres. BMC Med Genomics. 2019; 12(1):123. doi:10.1186/s12920-019-0561-0

Senol-Cosar O, Schmidt RJ, Qian E, Hoskinson D, Mason-Suares H, Funke B, et al. 6. Considerations for clinical curation, classification, and reporting of low-penetrance and low effect size variants associated with disease risk. Genet Med. 2019; 21(12):2765-73. doi:10.1038/s41436-019-0560-8

Otto PA, Horimoto ARVR. Penetrance rate estimation in autosomal dominant 7. conditions. Genet Mol Biol. 2012; 35(3):583-8. doi:10.1590/S1415-47572012005000051

Minikel EV, Vallabh SM, Lek M, Estrada K, Samocha KE, Sathirapongsasuti JF, et al. 8. Quantifying prion disease penetrance using large population control cohorts. Sci Transl Med. 2016; 8(322):322ra9. doi:10.1126/scitranslmed.aad5169

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. doi:10.1016/j.ajhg.2018.12.015 10. Chiò A, Battistini S, Calvo A, Caponnetto C, Conforti FL, Corbo M, et al. Genetic

counselling in ALS: facts, uncertainties and clinical suggestions. J Neurol Neurosurg Psychiatry. 2014; 85(5):478. doi:10.1136/jnnp-2013-305546

Iacoangeli A, Al Khleifat A, Sproviero W, Shatunov A, Jones AR, Opie-Martin S, et al. 11. ALSgeneScanner: a pipeline for the analysis and interpretation of DNA sequencing data of ALS patients. Amyotroph Lateral Scler Frontotemporal Degener. 2019; 20(3-4):207-15. doi:10.1080/21678421.2018.1562553

12. Abel O, Powell JF, Andersen PM, Al-Chalabi A. ALSoD: A user-friendly online bioinformatics tool for amyotrophic lateral sclerosis genetics. Hum Mutat. 2012; 33(9):1345-51. doi:10.1002/humu.22157

Al-Chalabi A, Lewis CM. Modelling the Effects of Penetrance and Family Size on Rates 13. of Sporadic and Familial Disease. Hum Hered. 2011; 71(4):281-8. doi:10.1159/000330167

14. Evans JDW, Girerd B, Montani D, Wang X-J, Galiè N, Austin ED, et al. BMPR2 mutations and survival in pulmonary arterial hypertension: an individual participant data meta-analysis. Lancet Respir Med. 2016; 4(2):129-37. doi:10.1016/S2213-2600(15)00544-5

Shatunov A, Al-Chalabi A. The genetic architecture of ALS. Neurobiol Dis. 2021; 15. 147:105156. doi:10.1016/j.nbd.2020.105156

16. Iacoangeli A, Al Khleifat A, Jones AR, Sproviero W, Shatunov A, Opie-Martin S, et al. C9orf72 intermediate expansions of 24–30 repeats are associated with ALS. Acta Neuropathol Commun. 2019; 7(1). doi:10.1186/s40478-019-0724-4

17. Hughes I, Hase T. Measurements and their uncertainties: a practical guide to modern error analysis. Oxford: Oxford University Press; 2010.

18. Zou Z-Y, Zhou Z-R, Che C-H, Liu C-Y, He R-L, Huang H-P. Genetic epidemiology of amyotrophic lateral sclerosis: a systematic review and meta-analysis. J Neurol Neurosurg Psychiatry 2017; 88:540-9. doi:10.1136/jnnp-2016-315018

Fertility rate, total (births per woman) [Internet]. 2020. Available from: 19. https://databank.worldbank.org/reports.aspx?source=2&series=SP.DYN.TFRT.IN

20. Shino MY, McGuire V, Van Den Eeden SK, Tanner CM, Popat R, Leimpeter A, et al. Familial aggregation of Parkinson's disease in a multiethnic community-based case-control study. Mov Disord. 2010; 25(15):2587-94. doi:10.1002/mds.23361

Elbaz A, Grigoletto F, Baldereschi M, Breteler MM, Manubens-Bertran JM, Lopez-21. Pousa S, et al. Familial aggregation of Parkinson's disease. Neurology. 1999; 52(9):1876. doi:10.1212/WNL.52.9.1876

22. Parkinson's UK. The Incidence and Prevalence of Parkinson's in the UK: Results from the Clinical Practice Research Datalink Reference Report. 2017. Available from: https://www.parkinsons.org.uk/professionals/resources/incidence-and-prevalenceparkinsons-uk-report

Healy DG, Falchi M, O'Sullivan SS, Bonifati V, Durr A, Bressman S, et al. Phenotype, 23. genotype, and worldwide genetic penetrance of LRRK2-associated Parkinson's disease: a case-control study. Lancet Neurol. 2008; 7(7):583-90. doi:10.1016/S1474-4422(08)70117-0

Thenappan T, Ryan JJ, Archer SL. Evolving epidemiology of pulmonary arterial 24. hypertension. Am J Respir Crit Care Med. 2012; 186(8):707-9. doi:10.1164/rccm.201207-1266ED

25. Aldred MA, Vijayakrishnan J, James V, Soubrier F, Gomez-Sanchez MA, Martensson G, et al. BMPR2 gene rearrangements account for a significant proportion of mutations in familial and idiopathic pulmonary arterial hypertension. Hum Mutat. 2006; 27(2):212-3. doi:10.1002/humu.9398

26. Byrne S, Walsh C, Lynch C, Bede P, Elamin M, Kenna K, et al. Rate of familial amyotrophic lateral sclerosis: a systematic review and meta-analysis. J Neurol Neurosurg Psychiatry. 2011; 82(6):623-7. doi:10.1136/jnnp.2010.224501

27. Byrne S, Heverin M, Elamin M, Bede P, Lynch C, Kenna K, et al. Aggregation of neurologic and neuropsychiatric disease in amyotrophic lateral sclerosis kindreds: A population-based case-control cohort study of familial and sporadic amyotrophic lateral sclerosis. Ann Neurol. 2013; 74(5):699-708. doi:10.1002/ana.23969

28. Marogianni C, Rikos D, Provatas A, Dadouli K, Ntellas P, Tsitsi P, et al. The role of C9orf72 in neurodegenerative disorders: a systematic review, an updated meta-analysis, and the creation of an online database. Neurobiol Aging. 2019:1.e-.e10. doi:10.1016/j.neurobiolaging.2019.04.012

Larkin EK, Newman JH, Austin ED, Hemnes AR, Wheeler L, Robbins IM, et al. 29. Longitudinal Analysis Casts Doubt on the Presence of Genetic Anticipation in Heritable Pulmonary Arterial Hypertension. Am J Respir Crit Care Med. 2012; 186(9):892-6. doi:10.1164/rccm.201205-0886OC

30. Andersen PM. Amyotrophic lateral sclerosis associated with mutations in the CuZn superoxide dismutase gene. Curr Neurol Neurosci Rep. 2006; 6(1):37-46. doi:10.1007/s11910-996-0008-9

31. Cudkowicz ME, McKenna-Yasek D, Sapp PE, Chin W, Geller B, Hayden DL, et al. Epidemiology of mutations in superoxide dismutase in amyotrophic lateral sclerosis. Ann Neurol. 1997; 41(2):210-21. doi:10.1002/ana.410410212

32. van Blitterswijk M, van Es MA, Hennekam EAM, Dooijes D, van Rheenen W, Medic J, et al. Evidence for an oligogenic basis of amyotrophic lateral sclerosis. Hum Mol Genet. 2012; 21(17):3776-84. doi:10.1093/hmg/dds199

33. Beck J, Poulter M, Hensman D, Rohrer JD, Mahoney CJ, Adamson G, et al. Large C9orf72 hexanucleotide repeat expansions are seen in multiple neurodegenerative syndromes and are more frequent than expected in the UK population. Am J Hum Genet. 2013; 92(3):345-53. doi:10.1016/j.ajhg.2013.01.011

Murphy NA, Arthur KC, Tienari PJ, Houlden H, Chio A, Traynor BJ. Age-related 34. penetrance of the C9orf72 repeat expansion. Sci Rep. 2017; 7(1):2116. doi:10.1038/s41598-017-02364-1

35. Coyle-Gilchrist ITS, Dick KM, Patterson K, Vázguez Rodríguez P, Wehmann E, Wilcox A, et al. Prevalence, characteristics, and survival of frontotemporal lobar degeneration syndromes. Neurology. 2016; 86(18):1736. doi:10.1212/WNL.00000000002638

Alonso A, Logroscino G, Jick SS, Hernán MA. Incidence and lifetime risk of motor 36. neuron disease in the United Kingdom: a population-based study. Eur J Neurol. 2009; 16(6):745-51. doi:10.1111/j.1468-1331.2009.02586.x

37. Turner MR, Al-Chalabi A, Chio A, Hardiman O, Kiernan MC, Rohrer JD, et al. Genetic screening in sporadic ALS and FTD. J Neurol Neurosurg Psychiatry. 2017; 88(12):1042-4. doi:10.1136/jnnp-2017-315995

Majounie E, Renton AE, Mok K, Dopper EGP, Waite A, Rollinson S, et al. Frequency of 38. the C9orf72 hexanucleotide repeat expansion in patients with amyotrophic lateral sclerosis and frontotemporal dementia: a cross-sectional study. Lancet Neurol. 2012; 11(4):323-30. doi:10.1016/s1474-4422(12)70043-1

39. Chiò A, Mazzini L, Alfonso S, Corrado L, Canosa A, Moglia C, et al. The multistep hypothesis of ALS revisited. Neurology. 2018; 91(7):e635.

doi:10.1212/WNL.000000000005996

Vajda A, McLaughlin RL, Heverin M, Thorpe O, Abrahams S, Al-Chalabi A, et al. 40. Genetic testing in ALS: A survey of current practices. Neurology. 2017; 88(10):991-9. doi:10.1212/wnl.000000000003686