

## Cascade screening following a polygenic risk score test: what is the risk of a relative conditional on a high score of a proband?

Shai Carmi<sup>1</sup>

<sup>1</sup> Braun School of Public Health and Community Medicine, The Hebrew University of Jerusalem, Jerusalem, Israel

shai.carmi@huji.ac.il

### Abstract

Polygenic risk scores (PRSs) for predicting disease risk have become increasingly accurate, leading to rising popularity of PRS tests. Consider an individual whose PRS test has placed him/her at the top  $q$ -quantile of genetic risk. Recently, Reid et al. (Circ Genom Precis Med. 2021;14:e003262) have investigated whether such a finding should motivate cascade screening in the proband's siblings. Specifically, using data from the UK biobank, Reid et al. computed the empirical probability of a sibling of the proband to also have a PRS at the top  $q$ -quantile. In this short note, I use the liability threshold model to compute this probability analytically (for either a sibling of the proband or for a more distant relative), showing excellent agreement with the empirical results of Reid et al., including that this probability is disease-independent. Further, I compute the probability of the relative of the proband to be affected, as a function of the quantile threshold  $q$ , the proportion of variance explained by the score, and the disease prevalence.

### Introduction

Polygenic risk scores show great promise for personalized disease risk prediction. A polygenic risk score (PRS) for a disease is a count of the number of risk alleles carried by an individual, with each allele weighted by its effect size. For complex diseases, a PRS represents the cumulative risk generated by thousands or more variants, each of a small effect<sup>1</sup>. PRSs were empirically shown to explain a substantial proportion of the variance in disease liability, and individuals at the highest PRS percentiles were shown to have risk that can be even  $\approx 3x$  higher compared to the population mean<sup>2-4</sup>. This suggests the feasibility of personalized prevention and/or intervention in these individuals, whose high risk may not always manifest as traditional risk factors<sup>5</sup>.

In clinical genetics, cascade screening refers to the testing of relatives of individuals who were either found to have a disease with a genetic component or were found to carry a disease mutation<sup>6</sup>. Cascade screening is important, as it allows high-specificity identification of individuals at high genetic risk. Cascade screening is currently limited to severe diseases and single mutations of large effects<sup>7</sup>. However, the development of increasingly accurate polygenic risk scores<sup>8</sup>, along with the increasing number of individuals receiving PRS results worldwide<sup>9</sup>, raise the prospects of performing cascade screening already upon the finding of an extreme PRS.

Recently, Reid et al. have investigated the question of cascade screening following a PRS finding<sup>10</sup>. They empirically studied the following setting. Suppose an individual has taken a PRS test that has placed

him/her at the top  $q$ -quantile of the PRS distribution (where  $q$  can be, e.g., 1%, 5%, 10%, etc.). What is the probability of a sibling of the proband to also have PRS at the top  $q$ -quantile? Using data on  $\approx 23,000$  sib-pairs from the UK biobank, Reid et al. have shown that the siblings of those probands have an increased probability to have a high PRS compared to the population average. Further, the probability was about the same across four diseases.

Despite these interesting results, several questions remain open. Specifically, it is unclear (1) whether the empirical results of Reid et al. (and particularly, the independence across diseases) are concordant with predictions of quantitative genetics theory; (2) what is the expected risk for diseases or threshold quantiles not empirically studied; (3) what is the expected risk for more distant relatives; and (4) what is the expected risk of the relative to become affected, rather than just have a PRS above a cutoff. The latter is important, as the cost-effectiveness of cascade screening depends on the actual disease risk rather than just the PRS. Here, I use the liability threshold model, an established model of disease risk in quantitative genetics, to address these questions.

## Methods

### Model

The liability threshold model (LTM) is a classic model in quantitative genetic theory that relates the risk of a disease to underlying genetic and non-genetic factors<sup>11,12</sup>. According to the LTM, a disease has an underlying continuous liability, distributed as a standard normal random variable. The liability can be written as  $y = g + \epsilon$ , where  $g \sim N(0, h^2)$  represents polygenic genetic factors (with variance equal to the heritability  $h^2$ ) and  $\epsilon \sim N(0, 1 - h^2)$  represents non-genetic (environmental) factors. An individual is affected whenever his/her liability exceeds a threshold. For a disease with prevalence  $K$ , the threshold is  $z_K$ , the upper- $K$  quantile of the standard normal distribution. (For example, for  $K = 0.01$ ,  $z_K = 2.33$ .) The LTM was found to fit well genetic data on complex diseases, and it is widely applied<sup>13-19</sup>.

In our setting, we do not know the precise value of the genetic factors influencing the liability. Instead, we have an estimate represented by the PRS. We thus write the liability as  $y = s + e$ , where  $s$  is the PRS and  $e$  is the residual liability, representing non-modeled genetic factors and non-genetic factors<sup>17,19,20</sup>. We assume that the PRS is normally distributed and has been standardized, such that  $s \sim N(0, r^2)$ , where  $r^2$  is the variance in liability explained by the PRS. Consequently,  $e \sim N(0, 1 - r^2)$ .

Consider next a pair of (full) siblings. Using standard quantitative genetic theory, it can be shown that that the PRSs of the two sibs can be written as  $s_1 = c + x_1$  for sib 1 and  $s_2 = c + x_2$  for sib 2. In these equations,  $c$  is a genetic component shared between the sibs, equal to the average maternal and paternal PRSs. Its distribution across the population is  $c \sim N(0, r^2/2)$ . Then,  $x_1 \sim x_2 \sim N(0, r^2/2)$  are two *independent* genetic components. For details on the derivation, see our previous publications<sup>21,22</sup>. For an intuitive explanation, the “segregation” variance, i.e., the variance of any polygenic component across siblings due to the randomness of meiosis, is known to be half the variance in the population<sup>23</sup>. Thus, given the parental PRSs, the PRS of each child has variance  $r^2/2$ .

We next define a threshold above which an individual is designated as having “high PRS”. We define the threshold as the upper- $q$  quantile of the PRS distribution. For example, if  $q = 0.01$ , an individual is considered to have high PRS if his/her PRS is at the top 1% of PRSs across the population. The PRS has

zero mean and variance  $r^2$  in the population, and thus, the value of the PRS at the threshold is  $z_q r$ , where  $z_q$  is the upper- $q$  quantile of a standard normal variable. We are told that sib 1 (the proband) has a high PRS, i.e.,  $s_1 > z_q r$ . We would like to compute the conditional probability that sib 2 either also has high PRS, or is affected.

Our calculations are similar to those of So et al.<sup>17</sup> and Do et al.<sup>19</sup>, who have considered the problem of predicting disease risk based on the PRS of an individual and/or a relative, along with the disease status of the relative. Here, we assume the disease status of the proband is unknown (e.g., for a late-onset disease).

### The probability that the sibling has high PRS

We would like to compute the probability  $P(s_2 > z_q r \mid s_1 > z_q r)$ . Using the definition of the conditional probability,

$$(1) P(s_2 > z_q r \mid s_1 > z_q r) = \frac{P(s_1 > z_q r \cap s_2 > z_q r)}{P(s_1 > z_q r)} = \frac{1}{q} P(s_1 > z_q r \cap s_2 > z_q r).$$

To compute the numerator, we condition on  $c$ . Given  $c$ , the scores of the two sibs are independent, i.e.,

$$(2) P(s_1 > z_q r \cap s_2 > z_q r \mid c) = P(s_1 > z_q r \mid c) P(s_2 > z_q r \mid c).$$

For  $i = 1, 2$ ,

$$(3) P(s_i > z_q r \mid c) = P(c + x_i > z_q r) = P(x_i > z_q r - c) = 1 - \Phi\left(\frac{z_q r - c}{r/\sqrt{2}}\right) = 1 - \Phi\left(\sqrt{2}z_q - \frac{c}{r/\sqrt{2}}\right).$$

In Eq. (3),  $\Phi(x)$  is the cumulative distribution of the standard normal variable, and we used the fact that  $x_i \sim N(0, r^2/2)$ . We can now compute the desired probability (Eq. (1)) by substituting Eq. (3) into Eq. (2), and integrating over all  $c$ . Recalling that  $c \sim N(0, r^2/2)$ ,

$$(4) P(s_2 > z_q r \mid s_1 > z_q r) = \frac{1}{q} \int_{-\infty}^{\infty} \frac{\phi\left(\frac{c}{r/\sqrt{2}}\right)}{r/\sqrt{2}} \left[1 - \Phi\left(\sqrt{2}z_q - \frac{c}{r/\sqrt{2}}\right)\right]^2 dc.$$

In Eq. (4),  $\phi(x)$  is the density of the standard normal variable. We now change variables,  $t = c/(r/\sqrt{2})$ , and obtain

$$(5) P(s_2 > z_q r \mid s_1 > z_q r) = \frac{1}{q} \int_{-\infty}^{\infty} \phi(t) \left[1 - \Phi(\sqrt{2}z_q - t)\right]^2 dt.$$

Eq. (5) is our final result for the probability of the sibling of the proband to have high PRS. Note that the probability does not depend on  $r^2$ , the variance explained by the PRS, and thus is disease- and PRS-independent.

### The probability that the sibling is affected

Denote by  $y_1$  and  $y_2$  the liabilities of the two sibs, and recall that sib 2 will be affected if  $y_2 > z_K$ , where  $z_K$  is the upper- $K$  quantile of the standard normal distribution and  $K$  is the prevalence. We would like to compute the following probability,

$$(6) P(y_2 > z_K \mid s_1 > z_q r) = \frac{P(y_2 > z_K \cap s_1 > z_q r)}{P(s_1 > z_q r)} = \frac{1}{q} P(y_2 > z_K \cap s_1 > z_q r).$$

As above, we write  $s_i = c + x_i$  for  $i = 1, 2$ , and condition on  $c$ . First, we note that

$$(7) P(y_2 > z_K \cap s_1 > z_q r \mid c) = P(y_2 > z_K \mid c) P(s_1 > z_q r \mid c),$$

because knowledge of the score of sib 1, given  $c$ , is not informative on the remaining liability components of sib 2. As in Eq. (3) above,

$$(8) P(s_1 > z_q r \mid c) = 1 - \Phi\left(\sqrt{2}z_q - \frac{c}{r/\sqrt{2}}\right).$$

Next,

$$(9) P(y_2 > z_K \mid c) = P(s_2 + e_2 > z_K \mid c) = P(c + x_2 + e_2 > z_K) = P(x_2 + e_2 > z_K - c) \equiv P(\tilde{e}_2 > z_K - c) = 1 - \Phi\left(\frac{z_K - c}{\sqrt{1 - r^2/2}}\right).$$

Above, we defined  $\tilde{e}_2 \equiv x_2 + e_2$ , and used the fact that  $x_2$  and  $e_2$  are independent normals, such that  $\tilde{e}_2$  is normal with zero mean and variance  $\text{Var}(\tilde{e}_2) = r^2/2 + (1 - r^2) = 1 - r^2/2$ . We can compute the desired probability (Eq. (6)) by substituting Eqs. (8) and (9) into Eq. (7), and integrating over all  $c$ ,

$$(10) P(y_2 > z_K \mid s_1 > z_q r) = \frac{1}{q} \int_{-\infty}^{\infty} \frac{\phi\left(\frac{c}{r/\sqrt{2}}\right)}{r/\sqrt{2}} \left[1 - \Phi\left(\sqrt{2}z_q - \frac{c}{r/\sqrt{2}}\right)\right] \left[1 - \Phi\left(\frac{z_K - c}{\sqrt{1 - r^2/2}}\right)\right] dc.$$

We can again change variables,  $t = c/(r/\sqrt{2})$ , and obtain

$$(11) P(y_2 > z_K \mid s_1 > z_q r) = \frac{1}{q} \int_{-\infty}^{\infty} \phi(t) \left[1 - \Phi(\sqrt{2}z_q - t)\right] \left[1 - \Phi\left(\frac{z_K - tr/\sqrt{2}}{\sqrt{1 - r^2/2}}\right)\right] dt.$$

This is our final expression for the probability of the sibling of the proband to be affected. Here, the probability depends on  $r^2$ , as well as on the prevalence  $K$ .

### The probability that an arbitrary relative has high PRS

We now turn from siblings into  $d$ -degree relatives. [Parents and children and full siblings are first degree relatives; half-sibs, grandparent and grandchild, and uncle and nephew are second-degree relatives; (full) first-cousins are third-degree relatives, and so on.] To compute the probability of the relative to be affected in this more general case, we take a different approach. Denote by  $s_1$  the PRS of the proband and by  $s_2$  the PRS of the relative. It is well known in quantitative genetics<sup>12,21</sup> that the distribution of  $(s_1, s_2)$  is multivariate normal, with the following parameters,

$$(12) (s_1, s_2) \sim MVN(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \boldsymbol{\mu} \equiv \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \boldsymbol{\Sigma} \equiv \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} r^2 & \frac{r^2}{2^d} \\ \frac{r^2}{2^d} & r^2 \end{pmatrix}.$$

The covariance term is  $\text{Cov}(s_1, s_2) = 2^{-d} \cdot r^2$  (i.e., the correlation is  $\rho = 2^{-d}$ ) because  $2^{-d}$  is the relatedness coefficient between  $d$ -degree relatives. As above, we would like to compute the probability  $P(s_2 > z_q r \mid s_1 > z_q r)$  that the relative has high PRS given that the proband has high PRS (top  $q$ -quantile).

Based on properties of multivariate normal distributions, we have

$$(13) s_2 \mid s_1 = s \sim N\left(\mu_2 + \frac{\sigma_2}{\sigma_1}\rho(s - \mu_1), (1 - \rho^2)\sigma_2^2\right) = N(\rho s, (1 - \rho^2)r^2) = N\left(\frac{s}{2^d}, \left(1 - \frac{1}{2^{2d}}\right)r^2\right).$$

The tail probability of  $s_2$  conditional on  $s_1$  is

$$(14) P(s_2 > z_q r \mid s_1 = s) = 1 - \Phi\left(\frac{z_q r - \frac{s}{2^d}}{\sqrt{\left(1 - \frac{1}{2^{2d}}\right)r^2}}\right) = 1 - \Phi\left(\frac{z_q r - \frac{s}{2^d}}{r\sqrt{1 - \frac{1}{2^{2d}}}}\right).$$

We can now compute the desired probability by integrating over  $s_1$ , recalling that  $s_1 \sim N(0, r^2)$ . Denote the probability density function of  $s_1$  as  $f_{s_1}(\cdot)$ . We have

$$(15) P(s_2 > z_q r \mid s_1 > z_q r) = \frac{1}{q} P(s_2 > z_q r \cap s_1 > z_q r) = \int_{z_q r}^{\infty} P(s_2 > z_q r \mid s_1 = s) f_{s_1}(s) ds = \frac{1}{q} \int_{z_q r}^{\infty} \left[1 - \Phi\left(\frac{z_q r - \frac{s}{2^d}}{r\sqrt{1 - \frac{1}{2^{2d}}}}\right)\right] \frac{\phi\left(\frac{s}{r}\right)}{r} ds.$$

We now change variables,  $t = s/r$ , and obtain

$$(16) P(s_2 > z_q r \mid s_1 > z_q r) = \frac{1}{q} \int_{z_q}^{\infty} \left[1 - \Phi\left(\frac{z_q r - \frac{tr}{2^d}}{r\sqrt{\left(1 - \frac{1}{2^{2d}}\right)}}\right)\right] \phi(t) dt = \frac{1}{q} \int_{z_q}^{\infty} \left[1 - \Phi\left(\frac{z_q - 2^{-d}t}{\sqrt{1 - 2^{-2d}}}\right)\right] \phi(t) dt = \frac{1}{q} \int_{z_q}^{\infty} \left[1 - \Phi\left(\frac{2^d z_q - t}{\sqrt{2^{2d} - 1}}\right)\right] \phi(t) dt.$$

Eq. (16) is our final result for the probability that the  $d$ -degree relative has high PRS given that the proband has high PRS. We validated numerically that for full siblings ( $d = 1$ ), Eq. (16) gives the same result as Eq. (5). Note that here too, the probability does not depend on  $r$  and hence is disease- and PRS-independent.

### The probability that an arbitrary relative is affected

Finally, we compute the probability that the relative of the proband is affected. As above, we denote by  $y_1$  and  $y_2$  the liabilities of the proband and relative, respectively, and recall that the relative will be affected if  $y_2 > z_K$ . As in Eq. (15) above,

$$(17) P(y_2 > z_K | s_1 > z_q r) = \frac{1}{q} P(y_2 > z_K \cap s_1 > z_q r) = \frac{1}{q} \int_{z_q r}^{\infty} P(y_2 > z_K | s_1 = s) f_{s_1}(s) ds.$$

Next,

$$(18) P(y_2 > z_K | s_1 = s) = \int_{-\infty}^{\infty} P(y_2 > z_K | s_1 = s, s_2 = s') P(s_2 = s' | s_1 = s) ds' = \int_{-\infty}^{\infty} P(y_2 > z_K | s_2 = s') P(s_2 = s' | s_1 = s) ds'.$$

The last step follows because given that we know the PRS of the relative, the total liability (and henceforth the disease status) of the relative no longer depend on the PRS on the proband. Recall that the liability is  $y = s + e$ , where  $e \sim N(0, 1 - r^2)$ . Thus,

$$(19) P(y_2 > z_K | s_2 = s') = P(s' + e > z_K) = P(e > z_K - s') = 1 - \Phi\left(\frac{z_K - s'}{\sqrt{1 - r^2}}\right).$$

Using Eq. (13) above,

$$(20) P(s_2 = s' | s_1 = s) = \frac{\phi\left(\frac{s' - \frac{s}{2^d}}{r\sqrt{1 - \frac{1}{2^{2d}}}}\right)}{r\sqrt{1 - \frac{1}{2^{2d}}}}.$$

Substituting Eqs. (19) and (20) into Eq. (18) we obtain

$$(21) P(y_2 > z_K | s_1 = s) = \int_{-\infty}^{\infty} \left[1 - \Phi\left(\frac{z_K - s'}{\sqrt{1 - r^2}}\right)\right] \frac{\phi\left(\frac{s' - \frac{s}{2^d}}{r\sqrt{1 - \frac{1}{2^{2d}}}}\right)}{r\sqrt{1 - \frac{1}{2^{2d}}}} ds' = \int_{-\infty}^{\infty} \left[1 - \Phi\left(\frac{z_K - r\sqrt{1 - 2^{-2d}}t'}{\sqrt{1 - r^2}}\right)\right] \phi\left(\frac{r\sqrt{1 - \frac{1}{2^{2d}}}t' - \frac{s}{2^d}}{r\sqrt{1 - \frac{1}{2^{2d}}}}\right) dt' = \int_{-\infty}^{\infty} \left[1 - \Phi\left(\frac{z_K - r\sqrt{1 - 2^{-2d}}t'}{\sqrt{1 - r^2}}\right)\right] \phi\left(t' - \frac{s}{r\sqrt{2^{2d} - 1}}\right) dt',$$

where we changed variables,  $t' = s' / \left(r\sqrt{1 - \frac{1}{2^{2d}}}\right)$ . We finally plug Eq. (21) into Eq. (17), recalling again that  $s_1 \sim N(0, r^2)$ . This gives

$$(22) P(y_2 > z_K | s_1 > z_q r) = \frac{1}{q} \int_{z_q r}^{\infty} \left\{ \int_{-\infty}^{\infty} \left[1 - \Phi\left(\frac{z_K - r\sqrt{1 - 2^{-2d}}t'}{\sqrt{1 - r^2}}\right)\right] \phi\left(t' - \frac{s}{r\sqrt{2^{2d} - 1}}\right) dt' \right\} \frac{\phi\left(\frac{s}{r}\right)}{r} ds.$$

Changing variables,  $t = s/r$ , we obtain

$$(23) P(y_2 > z_K | s_1 > z_q r) = \frac{1}{q} \int_{z_q}^{\infty} \left\{ \int_{-\infty}^{\infty} \left[ 1 - \Phi \left( \frac{z_K - r \sqrt{1-2^{-2d}t'}}{\sqrt{1-r^2}} \right) \right] \phi \left( t' - \frac{tr}{r\sqrt{2^{2d}-1}} \right) dt' \right\} \phi(t) dt =$$

$$\frac{1}{q} \int_{z_q}^{\infty} \left\{ \int_{-\infty}^{\infty} \left[ 1 - \Phi \left( \frac{z_K - r \sqrt{1-2^{-2d}t'}}{\sqrt{1-r^2}} \right) \right] \phi \left( t' - \frac{t}{\sqrt{2^{2d}-1}} \right) dt' \right\} \phi(t) dt.$$

Eq. (23) is our final equation for the probability of the relative to be affected. For the case of full siblings ( $d = 1$ ), we validated numerically that Eq. (23) gives the same result as Eq. (11). As for siblings, the probability of disease depends  $r^2$  and  $K$ .

## R implementation

For siblings, we solved the integrals in Eqs. (5) and (11) numerically using the function `integrate` in R. Our code is as follows.

```
risk_sib_high_prs = function(q) {
  integrand = function(t)
    return(dnorm(t)*pnorm(qnorm(1-q)*sqrt(2)-t, lower.tail=F)^2 / q)
  return(integrate(integrand, -Inf, Inf)$value)
}

risk_sib_affected = function(q,K,r) {
  zq = qnorm(1-q)
  zK = qnorm(1-K)
  integrand = function(t) {
    arg1 = zq*sqrt(2)-t
    arg2 = (zK-t*r/sqrt(2)) / (sqrt(1-r^2/2))
    return(dnorm(t)*pnorm(arg1,lower.tail=F)*pnorm(arg2,lower.tail=F) / q)
  }
  return(integrate(integrand, -Inf, Inf)$value)
}
```

For  $d$ -degree relatives, we again solved the integrals in Eqs. (16) and (23) numerically in R. Our code is as follows.

```
risk_rel_high_prs = function(q,d) {
  zq = qnorm(1-q)
  e = 2^d
  integrand = function(t)
    return(dnorm(t)*pnorm((2^d*zq-t)/sqrt(2^(2*d)-1), lower.tail=F) / q)
  return(integrate(integrand, zq, Inf)$value)
}

risk_rel_affected = function(q,K,r2,d)
{
  r = sqrt(r2)
  zq = qnorm(1-q)
  zK = qnorm(1-K)
  integrand_inner = function(tp, t) {
    arg1 = (zK - r*sqrt(1-2^(-2*d)))*tp/sqrt(1-r^2)
    arg2 = tp - t/sqrt(2^(2*d)-1)
    return(pnorm(arg1,lower.tail=F)*dnorm(arg2))
  }
}
```

```
integrand_outer = function(ts) {  
  y = numeric(length(ts))  
  for (i in seq_along(ts))  
  {  
    t = ts[i]  
    inner = integrate(integrand_inner, -Inf, Inf, t)$value  
    y[i] = dnorm(t)*inner / q  
  }  
  return(y)  
}  
return(integrate(integrand_outer, zq, Inf)$value)  
}
```

## Results and discussion

Reid et al.<sup>10</sup> have first studied the correlation between the PRSs of relatives. They found that the correlation between the PRSs of siblings and second-degree relatives was  $\approx 0.5$  and  $\approx 0.25$ , respectively. These correlations are naturally expected: based on standard quantitative genetic theory, the correlation between the genetic values of relatives is equal to their coefficient of relatedness<sup>12,24</sup>.

Next, Reid et al. have computed the empirical probability (across sib pairs in the UK biobank), that, given that the proband has PRS at the top  $q$  quantile, a sibling of the proband also has PRS at the top  $q$  quantile. Specifically, they considered four diseases (atrial fibrillation, coronary artery disease, diabetes, and severe obesity) and four values of  $q$  (1%, 5%, 10%, and 20%).

We used the liability threshold model to derive an analytical expression for the probability of the sibling of the proband to have PRS at the top  $q$  quantile (Methods; Eq. (5)). We compare our theoretical predictions to the empirical observations of Reid et al. in Figure 1A, showing excellent agreement. Our theory also predicts that the risk of the sibling is independent of the disease and the accuracy of the PRS, as empirically observed by Reid et al. Our figure also provides the expected risk for other values of  $q$  in the range 0-25%.

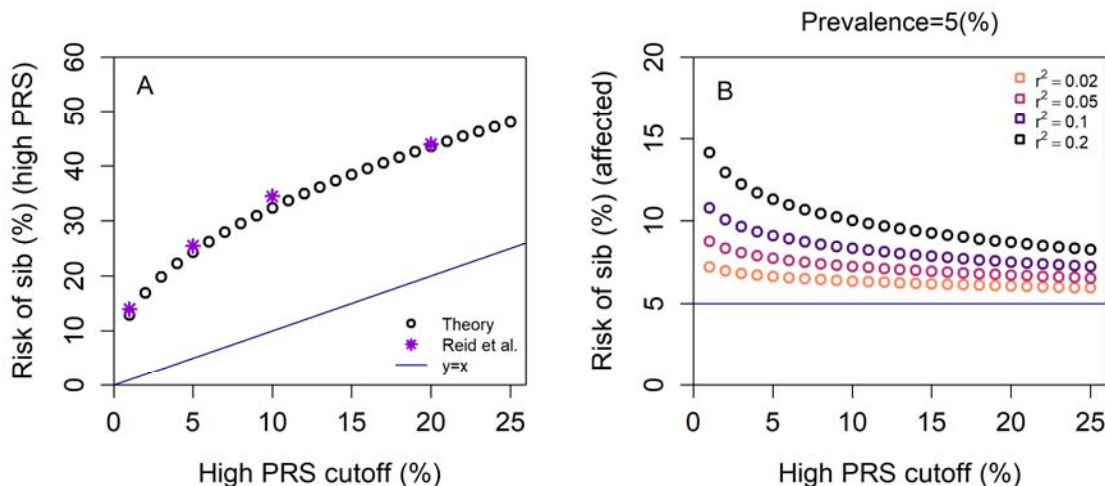
For studying the cost-effectiveness of cascade screening, it is necessary to estimate the risk of the sibling in the case screening is *not* applied. To this end, we need to compute the risk of the sibling to be affected, conditional on the proband having a high PRS. We used the liability threshold model to derive an analytical expression for this probability (Methods; Eq. (11)). We plot the results in Figure 1B, for a representative value of the prevalence (5%), and for four values of the proportion of variance explained by the PRS. As expected, the risk of the sibling is always higher than compared to a random individual from the population, and the risk increases with increasing accuracy of the score, and with a higher PRS of the proband (i.e., a smaller percentile used to define the top of the PRS distribution).

Full siblings are first degree relatives. To estimate the utility of cascade screening for more remote relatives of the proband, we extended our theory to compute the risk of a  $d$ -degree relative of the proband to have either high PRS (Eq. (16)) or to be affected (Eq. (23); Methods). While the resulting expressions are more complex, they are amenable to numerical evaluation. In Figure 2, we plot the risk of the relative to have high PRS (panel A) or to be affected (panel B) for  $d = 1, 2, 3, 4, 5$  ( $d = 3$

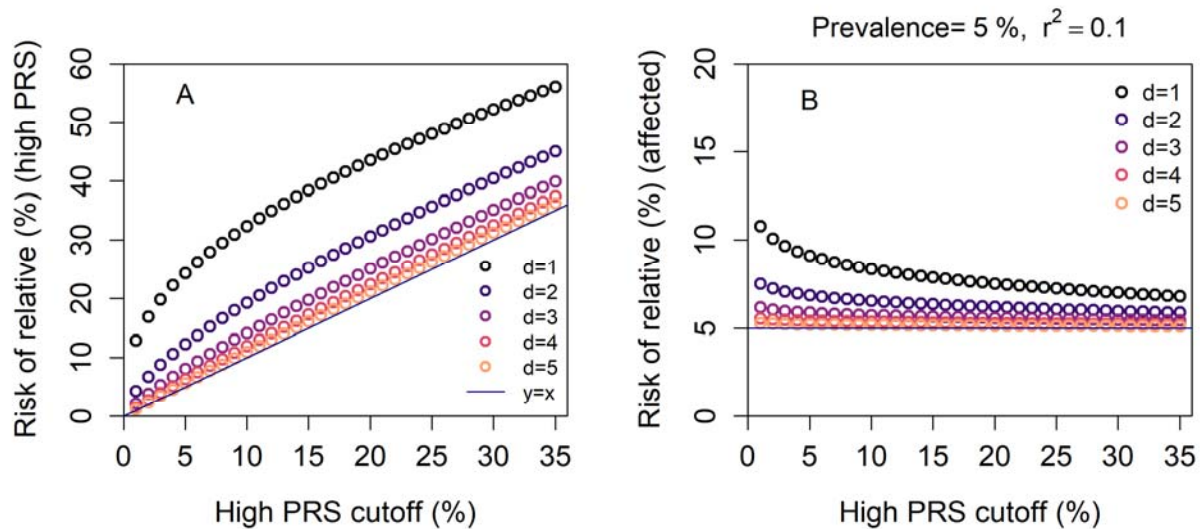


corresponds to first cousins and  $d = 5$  to second cousins). As expected, as  $d$  increases, the risk of the relative to have high PRS or to be affected approaches that of the general population.

Our simple R code will allow researchers to substitute any value for the high-PRS threshold, the accuracy of the score, the disease prevalence, and the degree of relatedness, in order to compute the expected outcomes of cascade screening in any setting of interest. Further conditioning on the disease status of the proband can be incorporated as in <sup>17,19</sup>. We expect our results to form a necessary building block for future studies of the cost-effectiveness of cascade screening following a PRS test.



**Figure 1. The expected risk of a sibling conditional on a proband having a PRS above a cutoff.** We assume that the proband is known to have a PRS at the top  $q$  quantile of the PRS distribution. The cutoff percentile varies along the x-axis. In (A), we plot the risk of the sib of the proband to have a high PRS, defined using the same cutoff. The diagonal blue line is  $y = x$ , which is the risk for an unrelated individual. The circles are the theoretical probabilities we derived based on the liability threshold model, obtained by numerically evaluating Eq. (5). The violet stars correspond to the empirical values (mean across diseases) obtained by Reid et al. In (B), we plot the risk of the sib to be affected. We assume prevalence of  $K = 5\%$ , and show results for multiple values of  $r^2$ , a measure of PRS accuracy equal to the proportion of variance in liability explained by the PRS (legend). The horizontal blue line represents the risk of an unrelated individual (equal to  $K$ ). The circles are the theoretical probabilities, obtained by numerically evaluating Eq. (11).



**Figure 2. The expected risk of a  $d$ -degree relative of the proband.** As in Figure 1, the x-axis represents the percentile cutoff used to define high PRS individuals. In (A), we plot the risk of the relative of the proband to have a high PRS, defined using the same cutoff. The diagonal blue line is  $y = x$ , which is the risk for an unrelated individual. The circles are the theoretical probabilities, obtained by numerically evaluating Eq. (16). Colors correspond to different degrees of relatedness  $d$  (legend). In (B), we plot the risk of the relative to be affected. We assume prevalence of  $K = 5\%$ , and show results for multiple values of  $d$  (legend). The horizontal blue line represents the risk of an unrelated individual (equal to  $K$ ). The circles are the theoretical probabilities (Eq. (23)).

## Competing interests

S. C. is paid consultant at MyHeritage.

## Bibliography

1. A. R. Martin, M. J. Daly, E. B. Robinson, S. E. Hyman, and B. M. Neale. Predicting Polygenic Risk of Psychiatric Disorders. *Biol Psychiatry* **86**, 97 (2019).
2. N. Mars, J. T. Koskela, P. Ripatti, T. T. J. Kiiskinen, A. S. Havulinna, J. V. Lindbohm *et al.* Polygenic and clinical risk scores and their impact on age at onset and prediction of cardiometabolic diseases and common cancers. *Nat Med* **26**, 549 (2020).
3. A. V. Khera, M. Chaffin, K. G. Aragam, M. E. Haas, C. Roselli, S. H. Choi *et al.* Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nat Genet* **50**, 1219 (2018).
4. L. Kachuri, R. E. Graff, K. Smith-Byrne, T. J. Meyers, S. R. Rashkin, E. Ziv *et al.* Pan-cancer analysis demonstrates that integrating polygenic risk scores with modifiable risk factors improves risk prediction. *Nat Commun* **11**, 6084 (2020).
5. L. Sun, L. Pennells, S. Kaptoge, C. P. Nelson, S. C. Ritchie, G. Abraham *et al.* Polygenic risk scores in cardiovascular risk prediction: A cohort study and modelling analyses. *PLoS Med* **18**, e1003498 (2021).
6. J. W. Knowles, D. J. Rader, and M. J. Khoury. Cascade Screening for Familial Hypercholesterolemia and the Use of Genetic Testing. *JAMA* **318**, 381 (2017).

7. P. J. Talmud, S. Shah, R. Whittall, M. Futema, P. Howard, J. A. Cooper *et al.* Use of low-density lipoprotein cholesterol gene score to distinguish patients with polygenic and monogenic familial hypercholesterolaemia: a case-control study. *Lancet* **381**, 1293 (2013).
8. S. A. Lambert, G. Abraham, and M. Inouye. Towards clinical utility of polygenic risk scores. *Hum Mol Genet* **28**, R133 (2019).
9. E. Widén, N. Junna, S. Ruotsalainen, I. Surakka, N. Mars, P. Ripatti *et al.* Communicating polygenic and non-genetic risk for atherosclerotic cardiovascular disease - An observational follow-up study. *medRxiv*, 2020.09.18.20197137 (2020).
10. N. J. Reid, D. G. Brockman, C. Elisabeth Leonard, R. Pelletier, and A. V. Khera. Concordance of a High Polygenic Score Among Relatives: Implications for Genetic Counseling and Cascade Screening. *Circ Genom Precis Med*, CIRCGEN120003262 (2021).
11. E. R. Dempster, and I. M. Lerner. Heritability of Threshold Characters. *Genetics* **35**, 212 (1950).
12. M. Lynch, and B. Walsh. *Genetics and analysis of quantitative traits*. (Sinauer Associates, 1998).
13. P. M. Visscher, and N. R. Wray. Concepts and Misconceptions about the Polygenic Additive Model Applied to Disease. *Hum Hered* **80**, 165 (2015).
14. N. R. Wray, and M. E. Goddard. Multi-locus models of genetic risk of disease. *Genome Med* **2**, 10 (2010).
15. S. H. Lee, N. R. Wray, M. E. Goddard, and P. M. Visscher. Estimating missing heritability for disease from genome-wide association studies. *Am J Hum Genet* **88**, 294 (2011).
16. O. Weissbrod, J. Flint, and S. Rosset. Estimating SNP-Based Heritability and Genetic Correlation in Case-Control Studies Directly and with Summary Statistics. *Am J Hum Genet* **103**, 89 (2018).
17. H. C. So, J. S. Kwan, S. S. Cherny, and P. C. Sham. Risk prediction of complex diseases from family history and known susceptibility loci, with applications for cancer screening. *Am J Hum Genet* **88**, 548 (2011).
18. D. S. Falconer. The inheritance of liability to certain diseases, estimated from the incidence among relatives. *Ann Hum Genet* **29**, 51 (1965).
19. C. B. Do, D. A. Hinds, U. Francke, and N. Eriksson. Comparison of family history and SNPs for predicting risk of complex disease. *PLoS Genet* **8**, e1002973 (2012).
20. S. H. Lee, M. E. Goddard, N. R. Wray, and P. M. Visscher. A better coefficient of determination for genetic profile analysis. *Genet Epidemiol* **36**, 214 (2012).
21. E. Karavani, O. Zuk, D. Zeevi, N. Barzilai, N. C. Stefanis, A. Hatzimanolis *et al.* Screening Human Embryos for Polygenic Traits Has Limited Utility. *Cell* **179**, 1424 (2019).
22. T. Lencz, D. Backenroth, E. Granot-Hershkovitz, A. Green, K. Gettler, J. H. Cho *et al.* Utility of polygenic embryo screening for disease depends on the selection strategy. *Elife* **10**, e64716 (2021).
23. N. R. Wray, K. E. Kemper, B. J. Hayes, M. E. Goddard, and P. M. Visscher. Complex Trait Prediction from Genome Data: Contrasting EBV in Livestock to PRS in Humans: Genomic Prediction. *Genetics* **211**, 1131 (2019).
24. P. M. Visscher, and M. E. Goddard. From R.A. Fisher's 1918 Paper to GWAS a Century Later. *Genetics* **211**, 1125 (2019).