


The quantitative genetics of the prevalence of infectious diseases: hidden genetic variation due to indirect genetic effects dominates heritable variation and response to selection

Piter Bijma,^{1,*} Andries D. Hulst,^{1,2} and Mart C. M. de Jong ²

¹Animal Breeding and Genomics, Wageningen University and Research, Wageningen 6708 PB, The Netherlands

²Quantitative Veterinary Epidemiology, Wageningen University and Research, Wageningen 6708 PB, The Netherlands

*Corresponding author: Email: piter.bijma@wur.nl

Abstract

Infectious diseases have profound effects on life, both in nature and agriculture. However, a quantitative genetic theory of the host population for the endemic prevalence of infectious diseases is almost entirely lacking. While several studies have demonstrated the relevance of transmission of infections for heritable variation and response to selection, current quantitative genetics ignores transmission. Thus, we lack concepts of breeding value and heritable variation for endemic prevalence, and poorly understand response of endemic prevalence to selection. Here, we integrate quantitative genetics and epidemiology, and propose a quantitative genetic theory for the basic reproduction number R_0 and for the endemic prevalence of an infection. We first identify the genetic factors that determine the prevalence. Subsequently, we investigate the population-level consequences of individual genetic variation, for both R_0 and the endemic prevalence. Next, we present expressions for the breeding value and heritable variation, for endemic prevalence and individual binary disease status, and show that these depend strongly on the prevalence. Results show that heritable variation for endemic prevalence is substantially greater than currently believed, and increases strongly when prevalence decreases, while heritability of disease status approaches zero. As a consequence, response of the endemic prevalence to selection for lower disease status accelerates considerably when prevalence decreases, in contrast to classical predictions. Finally, we show that most heritable variation for the endemic prevalence is hidden in indirect genetic effects, suggesting a key role for kin-group selection in the evolutionary history of current populations and for genetic improvement in animals and plants.

Keywords: quantitative genetics; infectious diseases; response to selection; indirect genetic effects; disease transmission; artificial selection; breeding programs; R_0

Introduction

Pathogens have profound effects on life on earth, both in nature and agriculture, and also directly on the human population (Schrag and Wiener 1995; Russel 2013). In nature, infectious pathogens are a major force shaping evolution of populations by natural selection, both in animals and plants (reviewed in Karlsson et al. 2014; Ebert and Fields 2020). In livestock, the annual cost of fighting and controlling epidemic and endemic infectious diseases is substantial, and much greater than the annual value of genetic improvement (Rushton 2009; Knap and Doeschl-Wilson 2020). Moreover, while antimicrobials have revolutionized medicine, the rapid appearance of resistant strains has resulted in a global health problem, both in the human population and in livestock (EFSA 2012; Thanner et al. 2016). Thus, there is an urgent need for additional methods and tools to combat infectious diseases. For livestock and plant production, artificial genetic selection of (host) populations for infectious disease traits may

provide such a tool. To quantify and optimize the potential benefits of such selection, however, we need to understand the quantitative genetics of infectious disease traits.

The integration of quantitative genetics and epidemiology for livestock populations was pioneered by Bishop and co-workers. They demonstrated unexpected effects, such as responses to selection for gastro-intestinal parasite infections clearly greater than expected from ordinary quantitative genetics (Bishop and Stear 1997, 1999, 2003). Bishop and co-workers also identified the basic reproduction number, R_0 , as a key parameter for genetic improvement and demonstrated the need for further integration of quantitative genetics and epidemiology (e.g., MacKenzie and Bishop 1999, 2001; Bishop and MacKenzie 2003; Nieuwhof et al. 2009; see also Doeschl-Wilson et al. 2021). These studies clearly show that classical quantitative genetic approaches do not predict response to genetic selection for disease traits, because they ignore the feed-back dynamics in the transmission of the

Received: July 13, 2021. Accepted: August 18, 2021

© The Author(s) 2021. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial reproduction and distribution of the work, in any medium, provided the original work is not altered or transformed in any way, and that the work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

infection. Several later studies have demonstrated the relevance of these transmission dynamics for heritable variation and response to selection in the host population (Lipschutz-Powell *et al.* 2012; Anche *et al.* 2014; Tsairidou *et al.* 2019; Hulst *et al.* 2021), mostly using stochastic simulation.

However, despite the findings of Bishop *et al.* (see references above) and the availability of well-established epidemiological theory (e.g., Diekmann *et al.* 2012), a quantitative genetic theory of the host population for the prevalence of infectious diseases is almost entirely lacking. The current theoretical framework of quantitative genetics and the approaches for genetic selection against infectious diseases in livestock and crops are largely based on the individual host response, ignoring transmission dynamics of the infection in the population. Moreover, we lack general expressions for the breeding value and genetic variance in key epidemiological parameters, in particular, for the basic reproduction number R_0 , even though such parameters may have a genetic basis.

Infections for which recovery does not confer any long-lasting immunity typically show endemic behavior, where the infection remains present in the population. For such infections, the endemic prevalence is defined as the expected fraction of the population that is infected. Because we lack a theoretical quantitative genetic framework for infectious diseases, we do not know which genetic effects of the host population determine the prevalence of an infectious disease, and have no concepts of breeding value and heritable variation for endemic prevalence. Hence, we do not understand the response of the endemic prevalence to genetic selection for disease traits at present. The main parameter determining the prevalence of endemic infections is the basic reproduction number R_0 , defined as the average number of individuals that gets infected by a typical infected individual in an otherwise noninfected population. In this article, we will propose a quantitative genetic framework for heritable variation and response to selection for R_0 and for the endemic prevalence of infectious diseases.

Individual phenotypes for infectious diseases are often recorded as the binary infection status of an individual, zero indicating noninfected and one indicating infected. The prevalence of an infection is then defined as the fraction of individuals that is infected, which is the fraction of individuals that has infection status $y = 1$. Because the average value of individual binary infection status is equal to the fraction of individuals infected, response to genetic selection in binary infection status is identical to response in prevalence, and *vice versa*. Binary infection status (0/1) typically shows low heritability, which suggests that response to selection is limited, also for prevalence (Bishop and Woolliams 2010; Bishop *et al.* 2012; Martin *et al.* 2018).

Geneticists have long realized that the categorical distribution of binary traits does not agree well with quantitative genetic models for polygenic traits, such as the infinitesimal model (Fisher 1919). For this reason, models have been developed that link an underlying normally distributed trait to the observed binary phenotype, such as the threshold model (Dempster and Lerner 1950; Gianola 1982) and the equivalent generalized linear mixed model with a probit link function (e.g., de Villemereuil *et al.* 2016). In such models, the underlying scale is interpreted as causal, and genetic parameters are assumed to represent “biological constants” on this scale. The genetic parameters on the observed scale, in contrast, depend on the mean of the trait, and thus change with the mean even when the change in allele frequencies at causal loci is infinitesimally small. In a landmark paper, Robertson (1950) showed that the observed-scale heritability of binary traits reaches a maximum at a prevalence of 0.5, and approaches zero when the prevalence is close to 0 or 1.

Hence, observed-scale heritability vanishes when artificial selection moves prevalence close to zero, hampering further genetic change.

Infectious disease status, however, differs fundamentally from binary phenotypes for noncommunicable traits, such as, say, heart failure. Because pathogens can be transmitted between host individuals, either directly or via the environment, the infection status of an individual depends on the status of other individuals in the population. This suggests that indirect genetic effects (IGEs) may play a role, which would fundamentally alter heritable variation and response to selection (Griffing 1967; Moore *et al.* 1997; Wolf *et al.* 1998; Bijma and Wade 2008; Bijma 2011). Results of simulation studies indeed suggest that selection response in the prevalence of infectious diseases may differ qualitatively from response in noncommunicable traits (Nieuwhof *et al.* 2009; Doeschl-Wilson *et al.* 2011; Anche *et al.* 2014; Hulst *et al.* 2021), and this has also been observed in an actual population (Heringstad *et al.* 2007). Results of Hulst *et al.* (2021), for example, show that genetic selection may result in the eradication of an infection via the mechanism of herd immunity, just like with vaccination (Fine 1993). This result contradicts predictions based on the observed-scale heritability for noncommunicable binary traits, where heritability vanishes when prevalence approaches zero (Robertson 1950).

While quantitative geneticists and breeders typically focus on individual disease status and (implicitly) interpret prevalence as an average of individual trait values, epidemiologists interpret the endemic prevalence of an infectious disease as the result of a population-level process of transmission of the infection (Kermack and McKendrick 1927; Keeling and Rohani 2011; Diekmann *et al.* 2012). In the latter perspective, both R_0 and the prevalence are emergent properties of a population, similar to the size of a termite colony or the number of prey caught by a hunting pack, rather than an average of individual trait values. Because such emergent traits do not belong to single individuals, we cannot apply the common partitioning of individual phenotypic values into individual additive genetic values (breeding values) and nonheritable residuals (“environment”). Nevertheless, the genetic effects that determine the response to selection in an emergent trait and the heritable variation for an emergent trait can be defined based on the so-called total heritable variation (Bijma 2011). The total heritable variation in a trait is based on the individual genetic effects on the level of the emergent trait, rather than on a decomposition of individual trait values into genetic and residual effects. This suggests we can develop a quantitative genetic theory for the endemic prevalence of infectious diseases by combining epidemiological theory with the theory of total heritable variation.

Here, we propose a quantitative genetic theory for the basic reproduction number R_0 and for the endemic prevalence of infectious diseases. We first identify the genetic factors that determine the prevalence of an infectious disease. Similar to the threshold model, we will assume an underlying additive infinitesimal model for those genetic factors. However, the link between the underlying additive scale and the observed endemic prevalence will be founded in epidemiological theory, with a key role for R_0 . Subsequently, we investigate the population-level consequences of genetic variation in individual disease traits for R_0 and for the endemic prevalence. Next, we move to the individual level, and derive expressions for the breeding value and heritable variation, for R_0 , endemic prevalence and individual binary infection status, and show how these parameters depend on the level of the endemic prevalence. Results will show that heritable

variation for endemic prevalence increases when prevalence approaches zero, while heritability of individual infection status goes to zero. Then we investigate response to selection against individual binary infection status (0/1), and show that response of prevalence to selection accelerates considerably when prevalence goes down. Finally, we partition the breeding value for prevalence into direct and IGE, and show that most of the heritable variation in the endemic prevalence of the infection is indirect, and thus hidden to classical genetic analysis and selection. We focus solely on the development of quantitative genetic theory, and do not consider the statistical estimation of the genetic effects underlying prevalence. Such methods have been developed elsewhere (Anacleto et al. 2015; Biemans et al. 2017; Pooley et al. 2020).

The theory we develop here applies to endemic microparasitic infections, i.e., where transmission depends just on whether an individual is infected or not, and where the infection is endemic in the local population (e.g., farm). It may also apply to endemic macroparasitic infection, such as coccidiosis or parasites, but we do not study this here. Endemic infections are of daily concern to farmers, and the very fact that they are endemic indicates that the existing management tools are insufficient. Thus, those are the infections likely to be targeted by breeding. Examples include mastitis, infectious claw disorders, respiratory infections in young animals (young replacement stock, meat calves, and fattening pigs), and fecal-oral transmitted infections causing gastro-intestinal diseases (diarrhea and so on.), but also several endemic and potentially zoonotic infections, such as *Salmonella* spp. and *Campylobacter jejuni* in poultry, Hepatitis E virus and MRSA in pigs, and Bovine Tb, *Leptospira*, *Brucella*, and Johne's disease in cows.

Theory and Results

The genetic factors that determine R_0 and the endemic prevalence

We consider an endemic infectious disease, where individuals can either be susceptible (i.e., in the noninfected state), denoted by S , or in the infected state, denoted by I . We use corresponding symbols in italics to denote the number of individuals with that status. Thus, with a total of N individuals in the population in which the endemic takes place, S denotes the number of susceptible individuals, I the number of infected individuals, and $S + I = N$ (see Table 1 for a notation key). We will assume that infected individuals are also infectious, and can thus infect others. When individuals recover they become susceptible again. This model is known as the SIS compartmental model (Hethcote 1989), and was first discussed by Weiss and Dishon (1971; In the Discussion, we will consider the validity of our results for other compartmental models).

The prevalence (P) of an endemic infection is defined as the fraction of the population infected (Diekmann et al. 2012),

$$P = \frac{I}{N}. \quad (1)$$

When individual infection status is coded in a binary fashion, using $y = 0$ for noninfected individuals and $y = 1$ for infected individuals, the prevalence is also equal to the average individual infection status in the population,

$$P = \bar{y}. \quad (2)$$

The prevalence of an infectious disease is determined by R_0 . The R_0 is defined as the average number of individuals that get

infected by a typical (i.e., average) infected individual in an otherwise noninfected population, and is a property of the population (Kermack and McKendrick 1927; Anderson and May 1979; Diekmann et al. 1990). When $R_0 > 1$, an average infected individual on average infects more than one new individual in an infection free population, and the infection can persist in the population.

The prevalence of an endemic infection reaches an equilibrium value, known as the endemic prevalence, when a single typical infected individual on average infects one other individual ($R = 1$; the endemic steady state). In a population where all individuals are the same, i.e., in the absence of genetic heterogeneity among host individuals, this occurs when the product of R_0 and the fraction of contact individuals that is susceptible is equal to one; $R_0(1 - P) = 1$. For example, when $R_0 = 3$, an infected individual could in principle infect three other individuals. However, when only one-third of its contact individuals is susceptible (i.e., not infected), meaning $1 - P = 1/3$, then the (average) total number of individuals that becomes infected by a single infected individual (the effective reproduction number, R) equals $3 \times 1/3 = 1$. Hence, when $1 - P = 1/3$, an infected individual is on average replaced by a single newly infected individual, so that an equilibrium occurs at $P = 1 - 1/3 = 2/3$. In the absence of heterogeneity, therefore, the endemic prevalence is given by (Weiss and Dishon 1971),

$$P = 1 - 1/R_0. \quad (3)$$

Throughout, we will use the symbol P to denote the endemic prevalence. The actual prevalence tends to fluctuate around the equilibrium value because of random perturbations and transient effects, for example when new animals replace some of the resident animals. Equation (3) is an approximation when there is variation among individuals, which is commonly referred to as "heterogeneity" in the epidemiological literature, and which will be addressed in the section on the impact of genetic variation on the endemic prevalence below.

Figure 1 illustrates the relationship between the endemic prevalence and R_0 . When R_0 is smaller than one the endemic prevalence is zero (the infection is not present in the long run), and Equation (3) does not apply. For large R_0 the endemic prevalence asymptotes to 1. This threshold phenomenon, i.e., $P = 0$ when $R_0 < 1$ and $P > 0$ when $R_0 \geq 1$, is exact also with heterogeneity (Diekmann et al. 1990). Note that the curve is steeper the closer R_0 is to 1. This pattern will have considerable consequences for the relationship between the heritable variation in the endemic prevalence and the level of the endemic prevalence, as will be shown in the section on individual genetic effects for the endemic prevalence below.

Because the endemic prevalence is determined by R_0 (Equation 3), the response of prevalence to selection (on any criterion), i.e., the genetic change in the endemic prevalence from one (host) generation to the next, follows from the genetic change in R_0 . Thus, to measure the value of an individual with respect to response to selection, we should base this measure on the genetic impact of the individual on R_0 . In other words, the definition of an individual breeding value for endemic prevalence should be based on R_0 . The next step, therefore, is to find the individual genetic factors underlying R_0 .

In the absence of variation among individuals (heterogeneity), R_0 is the product of the transmission rate parameter (β) and the mean duration of the infectious period ($1/\alpha$; Kermack and McKendrick 1927; Diekmann et al. 1990),

Table 1 Notation key

Symbol	Meaning
N	Total number of individuals in the population in which the endemic takes place; $N = I + S$
I	Number of infected individuals in the population
S	Number of susceptible (i.e., noninfected) individuals in the population
P	Prevalence in the endemic equilibrium; $P = I/N$; $P = \bar{y}$
R_0	Basic reproduction number
R_P	Response of prevalence to selection, i.e., change in prevalence per generation
$\mathcal{R}_{0,i}$	R_0 -like quantity that determines the prevalence for type i (see Equation 20)
A	Breeding value
A_{ly}, A_{lp}, A_{lx}	Breeding value for the logarithm of susceptibility, infectivity and recovery rate
A_{lR_0}, A_{R_0}	Breeding value for the logarithm of R_0 ; breeding value for R_0
A_y, A_P	Breeding value for individual binary disease status; breeding value for prevalence
G_{R_0}	Genotypic value for R_0
G_y, G_P	Genotypic value for individual disease status; genotypic value for prevalence.
T_P^2	Ratio of variance in breeding value for prevalence over phenotypic variance in y
c	Effective contact rate. Without heterogeneity $R_0 = c$.
h_y^2	Heritability of individual binary disease status
y	Individual binary infection status; infected: $y = 1$; noninfected: $y = 0$.
i, j	Subscript denoting an individual.
α	Recovery rate (relative to a value of 1)
β_{ij}	Transmission rate parameter from individual j to individual i .
γ	Susceptibility (relative to a value of 1)
ϕ	Infectivity (per unit of time, relative to a value of 1)
$\bar{\phi}_{inf}$	Mean infectivity of the infected individuals in the endemic equilibrium
ϕ	Life time infectivity (relative to a value of 1)
ϕ_{typ}	ϕ for the typical infected individual in an otherwise noninfected population
i	Intensity of selection; selection differential expressed in SD units
$\rho_{A_y, \bar{y}}$	Accuracy of mass selection, correlation of A_y and \bar{y} in the selection candidates
$\sigma_{A_{ly}}^2$	Variance of A_{ly} among individuals; analogous for $\sigma_{A_{lp}}^2$ and $\sigma_{A_{lx}}^2$
$\sigma_{A_{ly}, A_{lp}}^2$	Covariance of A_{ly} and A_{lp} ; analogous for $\sigma_{A_{ly}, A_{lx}}^2$ and $\sigma_{A_{lp}, A_{lx}}^2$
$\sigma_{A_{R_0}}^2$	Variance of the breeding values for the logarithm of R_0
$\sigma_{A_P}^2$	Additive genetic variance for endemic prevalence
$\sigma_{A_y}^2$	Additive genetic variance in individual binary infection status
$\sigma_{A_{P_D}}, \sigma_{A_{P_I}}^2$	Direct and indirect additive genetic variance for endemic prevalence, respectively
$\sigma_{A_{P_D}, A_{P_I}}$	Direct-indirect additive genetic covariance for endemic prevalence

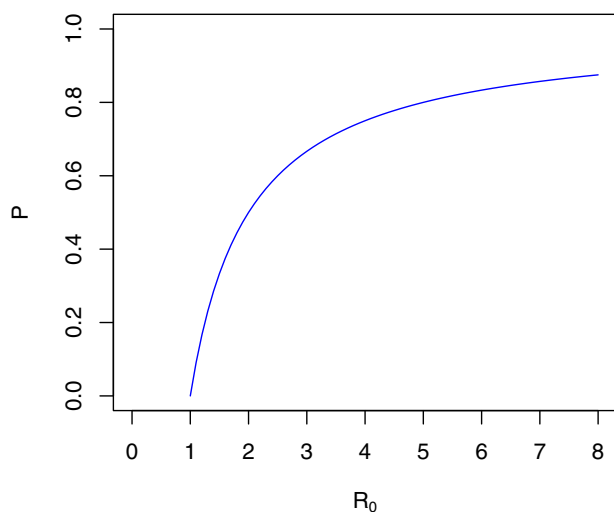


Figure 1 The relationship between the endemic prevalence (P) and the basic reproduction number (R_0) for a homogeneous population (from Equation 3).

$$R_0 = \beta/\alpha \quad (4)$$

where α is the recovery rate parameter. The β is the average number of individuals infected per unit of time by a single infected individual when all its contact individuals are susceptible, and α is the probability per unit of time for an infected individual to recover. With heterogeneity, Equation (4) is an approximation

(Diekmann et al. 1990); the effect of heterogeneity on R_0 will be addressed below.

With heterogeneity, the transmission rate parameter may vary between pairings (contacts) of individuals. The transmission rate parameter between infectious individual j and susceptible individual i may be modelled as the product of an overall effective contact rate (c) for the population, and two individual quantitative genetic traits: the susceptibility (γ) of recipient individual i and the infectivity (ϕ) of donor individual j (e.g., DeJong et al. 1996; Lipschutz-Powell et al. 2014; Anacleto et al. 2015; Biemans et al. 2017),

$$\beta_{ij} = c\gamma_i\phi_j \quad (5)$$

where β_{ij} refers to transmission from individual j to i , and may differ from β_{ji} . The overall effective contact rate c is the transmission rate parameter β_{ij} for the average pair of individuals ij , for which $\gamma_i = \phi_j = 1$. The susceptibility γ_i is the propensity of the noninfected individual i to become infected, expressed relative to a value of 1. Analogously, the infectivity ϕ_j the propensity of an infected individual j to infect another individual, expressed relative to value of 1. Note that i and j are distinct individuals in Equation (5), so that γ_i and ϕ_j are independent when individuals i and j are genetically unrelated.

Equations (3) through (5) show that the factors underlying the endemic prevalence of an infection are the contact rate c , the susceptibility, γ , the infectivity, ϕ , and the recovery rate α . We define c as a fixed parameter for the population (or, for example, for a sex, herd or age class combination), whereas γ , ϕ , and α are

quantitative traits that may show random variation among individuals. Note, while the actual contact rate may vary among individuals, it is convenient to include such variation in the individual susceptibility and infectivity traits. Thus, we also assume that all the individuals are mixing randomly within the (local) population, such as a herd. In principle, β_{ij} might also depend on the specific combination of i and j , so that we cannot fully separate β_{ij} into a product of components due to i and j . However, in a quantitative genetic perspective, such a combination effect represents interactions between genes in distinct individuals (say “between-individual epistasis”), which does not contribute to the heritable variation, and which we will therefore ignore. In epidemiological terminology, we assume separable mixing (Diekmann et al. 1990). Hence, conceptually we define c as the average effective contact rate for the population, while variation in contact rate among individuals is included in γ and ϕ . Moreover, to define the scale of Equations (4) and (5), it is convenient to include the scale in c , and to express γ , ϕ , and α relative to a value of 1. Hence, with this parameterization, the c is on the scale of R_0 , and R_0 and c are identical in the absence of heterogeneity. With heterogeneity, however, R_0 may deviate from c (see below).

Genetic models for susceptibility, infectivity, recovery, and R_0

Genetic variation is potentially present in susceptibility, infectivity, and the recovery rate. In this section, we propose a genetic model for these traits, which subsequently leads to a genetic model for R_0 .

We assume that susceptibility, infectivity, and the recovery rate are affected by a large number of loci, each of small effect, so that genetic effects approximately follow a normal distribution. However, as γ , ϕ , and α represent rates, i.e., probabilities per unit of time, their values are strictly positive. Moreover, in the expressions for the epidemiological parameters of the previous section (Equations 4 and 5), all these parameters appear in products. For these reasons, following Anacleto et al. (2015), we define normally distributed additive genetic effects for the logarithm of these rates, so that effects are multiplicative on the actual scale, and the rates themselves follow a log-normal distribution,

$$\gamma_i = e^{A_{\gamma,i}} \quad (6a)$$

$$\phi_i = e^{A_{\phi,i}} \quad (6b)$$

$$\alpha_i = e^{A_{\alpha,i}} \quad (6c)$$

where $A_{l,i}$ denotes the Normally distributed additive genetic value (breeding value) for the logarithm of the corresponding rate for individual i , and has a mean of zero,

$$\begin{bmatrix} A_{\gamma} \\ A_{\phi} \\ A_{\alpha} \end{bmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{A_{\gamma}}^2 & \sigma_{A_{\gamma}A_{\phi}} & \sigma_{A_{\gamma}A_{\alpha}} \\ \sigma_{A_{\gamma}A_{\phi}} & \sigma_{A_{\phi}}^2 & \sigma_{A_{\phi}A_{\alpha}} \\ \sigma_{A_{\gamma}A_{\alpha}} & \sigma_{A_{\phi}A_{\alpha}} & \sigma_{A_{\alpha}}^2 \end{pmatrix} \right). \quad (7)$$

Throughout, we use subscript l to denote the natural logarithm. Thus, the breeding values for $\log(\gamma)$, $\log(\phi)$, and $\log(\alpha)$ follow a multivariate normal distribution, as common in quantitative genetics. Moreover, for the average individual the $A_l = 0$, so that its rates are equal to one ($\gamma = \phi = \alpha = 1$). Hence, those rates should be interpreted relative to a value of 1. An individual with $\gamma = 2$, for example, is twice as susceptible as the average individual. Also note that, by defining breeding values to have a mean of zero, we put the mean into the contact rate c . Hence, in the following, c will refer to the model where the mean breeding

value on the log scale is equal to zero (Equation 7, see also Discussion).

The breeding values on the log-scale can approximately be interpreted as a relative change of the corresponding rate. For example, since $e^{0.1} \approx 1.1$, an A_{γ} of 0.1 corresponds approximately to a 10% greater than average susceptibility ($\gamma \approx 1.1$). Similarly, an A_{γ} of -0.1 corresponds approximately to a 10% smaller than average susceptibility ($\gamma \approx 0.9$). Realistic values for the genetic variances on the log-scale are probably smaller than $\sim 0.5^2$ (Hulst et al. 2021). For example, with $\sigma_{A_{\gamma}}^2 = 0.5^2$, the 10% least susceptible individuals have $\bar{\gamma} = e^{-0.88} = 0.42$, while the 10% most susceptible individuals have $\bar{\gamma} = e^{0.88} = 2.40$. Thus, the average susceptibilities of these top and bottom 10% of individuals differ by a factor of 5.7, which is substantial. Therefore, we will consider additive genetic variances on the log-scale no greater than 0.5^2 . With a prevalence of 0.3, this value corresponds to an observed-scale heritability of individual binary infection status of about 0.05 (Hulst et al. 2021).

Genotypic value and breeding value for R_0

Based on Equations (4) and (5), we may define an individual genotypic value for R_0 ,

$$G_{R_0,i} = c\gamma_i\phi_i/\alpha_i. \quad (8)$$

In contrast to the pair-wise transmission rate parameter β_{ij} in Equation (5), an individual's genotypic value for R_0 is entirely a function of its own rates, as can be seen from the index i on all elements of Equation (8). This is because $G_{R_0,i}$ refers to the genetic effects that originate from the individual, rather than to those that affect its trait value. As a consequence these rates may be correlated, as defined in Equation (7) above. Hence, $G_{R_0,i}$ represents a total genotypic value (Bijma et al. 2007; Bijma 2011). We focus on the total genotypic value, because our ultimate interest is in response to selection. In the next section of this manuscript, we will show that R_0 is indeed the simple population average of G_{R_0} .

From Equations (6a, b & c) and (8),

$$\begin{aligned} G_{R_0,i} &= c e^{A_{\gamma,i}} e^{A_{\phi,i}} / e^{A_{\alpha,i}} \\ &= e^{\ln(c) + A_{\gamma,i} + A_{\phi,i} - A_{\alpha,i}} \\ &= e^{\ln(c) + A_{R_0,i}}, \end{aligned} \quad (9)$$

where $A_{R_0,i}$ is a normally distributed additive genetic effect (breeding value) for the logarithm of R_0 ,

$$A_{R_0,i} = A_{\gamma,i} + A_{\phi,i} - A_{\alpha,i}, \quad (10a)$$

$$A_{R_0} \sim N(0, \sigma_{A_{R_0}}^2), \quad (10b)$$

$$\sigma_{A_{R_0}}^2 = \sigma_{A_{\gamma}}^2 + 2\sigma_{A_{\gamma}A_{\phi}} - 2\sigma_{A_{\gamma}A_{\alpha}} + \sigma_{A_{\phi}}^2 - 2\sigma_{A_{\phi}A_{\alpha}} + \sigma_{A_{\alpha}}^2. \quad (10c)$$

Hence, our model of the genotypic value for R_0 is additive with normally distributed effects on the log-scale. Thus, the genotypic value for R_0 , as defined in Equations (8) and (9), follows a log-normal distribution,

$$G_{R_0} \sim \log N(\mu = \ln(c), \sigma^2 = \sigma_{A_{R_0}}^2). \quad (11)$$

The genotypic value for R_0 for the average individual, which has $A_{R_0} = 0$, is equal to the contact rate, c . Hence, the genotypic value is defined here including its average, it is not expressed as a deviation from the mean. Moreover, we refer to G_{R_0} as a genotypic

value, rather than a breeding value, because the $e^{A_{IR_0}}$ in Equation (9) is a nonlinear function, so that G_{R_0} will show some nonadditive genetic variance, even though A_{IR_0} is additive.

The log-normal distribution of G_{R_0} agrees with the infinitesimal model and the strictly positive values for R_0 (Anacleto et al. 2015), and is also convenient because the mean and variance of G_{R_0} follow from the known properties of the log-normal distribution,

$$E(G_{R_0}) = c e^{\frac{1}{2}\sigma_{A_{IR_0}}^2} \quad (12)$$

$$\text{var}(G_{R_0}) = c^2(e^{2\sigma_{A_{IR_0}}^2} - e^{\sigma_{A_{IR_0}}^2}). \quad (13)$$

Equations (12) and (10c) show that genetic (co)variation in susceptibility, infectivity and/or recovery, and thus in the breeding value for the logarithm of R_0 , leads to an increase in the mean genotypic value for R_0 . For example, for $\sigma_{A_{IR_0}}^2 = 0.5^2$, $E(G_{R_0}) \approx 1.13c$. While this 13% increase in G_{R_0} may suggest limited impact of heterogeneity, a 13% increase in R_0 has a considerable impact on the endemic prevalence when R_0 is close to one (Figure 1).

Equations (12) and (13) show that a log-normal distribution for susceptibility, infectivity and recovery results in a positive mean-variance relationship for G_{R_0} . Figure 2 illustrates this relationship, for $\sigma_{A_{IR_0}}^2 = 0.3^2$ and genetic variation in susceptibility only. The x-axis shows the contact rate, which is equal to the genotypic value for R_0 of the average individual in the population (i.e., an individual with $\gamma_i = \varphi_i = \alpha_i = 1$). Hence, the x-axis reflects the level of R_0 . The small circle represents a population with a prevalence of ~ 0.33 , for which observed-scale heritability of binary infection status is ~ 0.02 (Hulst et al. 2021). For that population, R_0 is ~ 1.5 , and the genetic standard deviation in R_0 is ~ 0.48 . Hence, despite the small observed-scale heritability, R_0 has considerable genetic variation and some individuals will have a genotypic value smaller than 1, which agrees with the findings of Hulst et al. (2021). In the context of artificial selection against infectious diseases, the positive mean-variance relationship resulting from our model may be interpreted as conservative, because it implies a reduction of the genetic variance in R_0 with continued selection for lower prevalence.

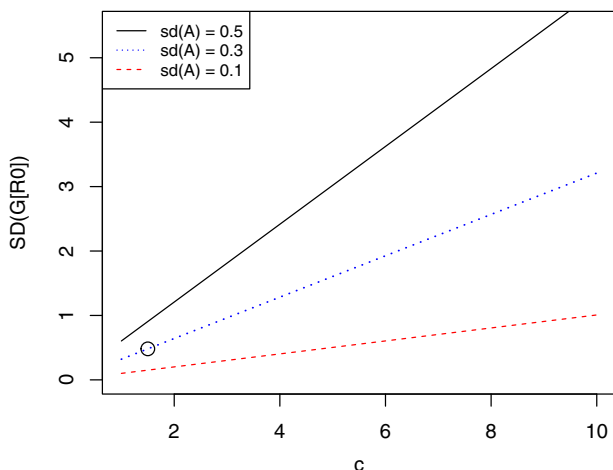


Figure 2 Genetic standard deviation in R_0 as a function of the level of R_0 (as approximated by the contact rate c here). From Equation (13), for three values of $\sigma_{A_{IR_0}}^2 = \sigma_{A_{\gamma}}^2 = 0.1^2, 0.3^2$, and 0.5^2 , and no variation in φ and α . The circle represents a population with a prevalence of ~ 0.33 , for which observed-scale heritability of binary infection status is ~ 0.02 (Hulst et al. 2021).

In summary, this section has presented a genetic model for susceptibility, infectivity and recovery, leading to expressions for the genotypic value and genetic variance in R_0 (Equations 8, 9, and 11). Note however, that we have not yet provided formal proof that the individual genotypic value for R_0 indeed predicts the actual R_0 (of the population so to say). In fact, the definition of G_{R_0} in Equation (8) is an educated guess based on the expression for R_0 in a homogeneous population (Equation 4). In epidemiology, however, R_0 is an emerging property of a population-level process of the transmission of an infection, rather than an average of individual (genotypic) values. Thus, it remains to be proven that the G_{R_0} defined in Equation (8) indeed predicts the R_0 of a genetically heterogeneous population. In the next two sections, therefore, we will focus on the population-level consequences of genetic heterogeneity, and investigate the impact of genetic variation on the level of R_0 and on the endemic prevalence.

The impact of genetic heterogeneity on R_0

R_0 is a key parameter for infectious diseases, because infections can persist in a population if and only if R_0 is greater than one (Kermack and McKendrick 1927; Diekmann et al. 1990). In other words, an endemic equilibrium can exist only when R_0 is greater than 1. Conversely, eradication of an infectious disease, either by vaccination or other measures such as genetic selection of the host population, requires that R_0 is reduced to a value smaller than one. Here, we address the consequences of genetic (co)variation in susceptibility, infectivity and recovery for the value of R_0 , and provide a proof that R_0 is indeed the simple population average of the individual genotypic values for R_0 , as defined in Equations (8), (9), and (11). Because our interest is in the impact of genetic heterogeneity on R_0 and in the genotypic value for R_0 , we consider genetic (co)variation only, disregarding environmental sources of (co)variation. Note that R_0 is strictly defined for the infection free state of the population (i.e., where the infected fraction is infinitesimally small). Hence, in this section we consider the infection free state, while the endemic equilibrium will be addressed in the next section.

Genetic (co)variation in susceptibility, infectivity and recovery has two consequences for R_0 . First, it increases the mean genotypic value for R_0 because the expectation of a log-normal variate increases with the variance on the log-scale. This effect is trivial; it follows directly from Equations (12) and (10c) and is not the main focus of this section. Second, as stated above, R_0 is the average number of individuals that gets infected by a typical infected individual in an otherwise noninfected (large) population (Kermack and McKendrick 1927; Diekmann et al. 1990). The expression for R_0 given in Equation (4) ignores the “typical” term in the definition of R_0 , and is therefore an approximation in case of heterogeneity (Diekmann et al. 1990). The focus of this section is on the consequences of heterogeneity for the properties of the typical infected individual, and thus for R_0 .

The properties of the “typical infected individual” will depend on the magnitude and nature of the heterogeneity among the individuals in the population, because the susceptibility and recovery determine which animals are infected, while the infectivity of those individuals may differ from the population average. In contrast to the conclusion of Springbett et al. (2003), therefore, genetic heterogeneity can affect R_0 (Diekmann et al. 1990, 2012). Suppose, for example, that individuals differ in both susceptibility and infectivity, and that susceptibility is positively correlated to infectivity. Because individuals with greater susceptibility are more likely to become infected, the typical infected individual

will have an above-average susceptibility. Moreover, because of the positive correlation with infectivity, this will also translate into an above average infectivity of the typical infected individual, leading to higher R_0 . Hence, variation among individuals together with a positive (negative) correlation between susceptibility and infectivity results in an increase (decrease) in R_0 (Diekmann et al. 1990). A similar argument holds for genetic covariation between recovery and infectivity. For this reason, R_0 in general deviates from the right-hand side of Equation (4) obtained using the averages of α , γ , and ϕ .

In Appendix A, we derive the relationship between R_0 and the genetic parameters for susceptibility, infectivity and recovery. The first step is the derivation of the lifetime infectivity of the typical infected individual. Lifetime infectivity is the total infectivity of an individual, aggregated over its average infectious period, and is the product of its infectivity per unit of time (ϕ , Equations 5 and 6) and the mean duration of its infectious period, $1/\alpha_i$,

$$\phi_i = \phi_i / \alpha_i. \quad (14)$$

Hence, $c\phi_i$ is the average total number of individuals that become infected by individual i over its entire infectious life time. By introducing life time infectivity, we summarize the infectivity (per unit of time, ϕ) and the recovery of an individual into a single variable (ϕ), which simplifies the analysis. Appendix A shows that the lifetime infectivity of the typical infected individual equals

$$\phi_{\text{typ}} = \bar{\phi} e^{\sigma_{A_{\gamma}, A_{\phi}}} \quad (15)$$

where $\bar{\phi}$ is the simple average of lifetime infectivity in the entire population, and $\sigma_{A_{\gamma}, A_{\phi}}$ the covariance between the breeding values for the logarithms of susceptibility and lifetime infectivity. Thus, Equation (15) is the average lifetime infectivity of infected individuals in the population at the (almost) infection free state, and can be interpreted as an average weighted for differences in susceptibility. It shows that the typical infected individual has an above (below) average life-time infectivity when the covariance between susceptibility and life-time infectivity is positive (negative), as argued verbally in the previous paragraph.

From the definition of R_0 and Equations (4), (5), and (15), it follows that

$$R_0 = c \bar{\gamma} \bar{\phi} e^{\sigma_{A_{\gamma}, A_{\phi}}} \quad (16)$$

where $\bar{\gamma}$ is the simple population average value of susceptibility. The last term of this expression shows that a positive covariance between susceptibility and life-time infectivity indeed increases R_0 . Note that this R_0 is the reproduction number in the large infection free population, as it is normally defined in epidemiology. Hence the distribution of susceptibility in the susceptible individuals remains equal to the population distribution, in contrast to the distribution of infectivity in the infected individuals (see Appendix A).

Equation (16) can be simplified by substituting the expression for $\bar{\gamma}$ and $\bar{\phi}$, which follow from the log-normal distribution,

$$\bar{\gamma} = e^{\frac{1}{2}\sigma_{A_{\gamma}}^2} \quad (17a)$$

$$\bar{\phi} = e^{\frac{1}{2}\sigma_{A_{\phi}}^2}. \quad (17b)$$

Substituting Equations (17a) and (17b) into Equation (16), and expressing the genetic variance of life-time infectivity in terms of infectivity per unit of time and recovery, reveals that R_0 is equal to the simple population average of the individual genotypic value for R_0 (see Appendix A and Equation 12),

$$R_0 = c e^{\frac{1}{2}\sigma_{A_{R_0}}^2} = E(G_{R_0}). \quad (18)$$

Thus, R_0 depends on the variance in the breeding value for the logarithm of R_0 , but is still equal to the mean genotypic value for R_0 . In other words, a positive covariance between susceptibility and life-time infectivity indeed increases R_0 , but this effect is fully captured by the effect of the variance in the breeding values for the logarithm of R_0 on the mean genotypic value for R_0 (the $e^{\frac{1}{2}\sigma_{A_{R_0}}^2}$ term in Equation 18). This result, therefore, provides formal proof that the genotypic value for R_0 , as defined in Equations (8)–(11), indeed represents the individual genetic value for R_0 .

Note that, while R_0 is equal to the simple average genotypic value for R_0 , it still differs from the product of the simple averages of the rates when susceptibility, infectivity and/or recovery are correlated; $R_0 \neq c \bar{\gamma} \bar{\phi} / \bar{\alpha}$. Moreover, R_0 may also differ from the simple $\bar{\beta} / \bar{\alpha}$ with heterogeneity. A numerical investigation of the $e^{\sigma_{A_{\gamma}, A_{\phi}}}$ term in Equation (16) shows that a correlation between susceptibility and life time infectivity may change R_0 by a maximum of about 25% for realistic levels of heterogeneity and log-normally distributed genetic effects. For example, for $\sigma_{A_{\gamma}}^2 = \sigma_{A_{\phi}}^2 = 0.5^2$, and a correlation $r_{A_{\gamma}, A_{\phi}} = 0.8$, R_0 is 22% greater than $c \bar{\gamma} \bar{\phi} / \bar{\alpha}$. For values of R_0 close to 1, this 22% may be the difference between absence of an infection vs. a significant endemic prevalence. Thus, a correlation between susceptibility, infectivity and/or recovery may have a meaningful impact on R_0 .

In summary, this section has shown that heterogeneity and a positive correlation between susceptibility and life-time infectivity lead to an increase of R_0 , and thus increase the probability that an infectious disease persists in the population. However, when genotypic values for R_0 follow a log-normal distribution, R_0 is still equal to the simple average of those genotypic values.

The impact of genetic variation on the endemic prevalence

In this section, we present an expression for the endemic prevalence in a population with genetic variation in susceptibility, infectivity and recovery, and also briefly investigate the quantitative effect of such variation for the endemic prevalence. Figure 1 and Equation (3) show the relationship between R_0 and the endemic prevalence for a homogeneous population. With variation among individuals, however, more susceptible individuals are more likely to be in the infected state in the endemic equilibrium. For this reason, the mean susceptibility of the remaining noninfected individuals will be lower than the population average susceptibility. This in turn translates into an endemic prevalence lower than expected based on R_0 [Equation 3; Springbett et al. 2003; Diekmann et al. 2012; Note, however, that the threshold value of $R_0 = 1$ remains, so that endemic prevalence is zero if and only if $R_0 \leq 1$, and in that sense the R_0 given in Equation (18) is exact]. Similar arguments can be used to show that prevalence depends on the variation in the recovery rate, and on the covariation of infectivity, susceptibility and recovery. Thus, Equation (3) is exact only in the absence of heterogeneity in these parameters.

The endemic prevalence in a heterogeneous population can be found by realizing that the prevalence must have reached an equilibrium value for each type of individual (Biemans *et al.* 2017; Aznar *et al.* 2018). Suppose, for example, that susceptibility, infectivity, and recovery would be governed by the same single bi-allelic locus in a diploid organism. Then, for the entire population to be in equilibrium, each of the three genotypic classes should be in equilibrium as well. In other words, the prevalence should have reached an equilibrium value within each genotypic class, but this value may differ among the three classes. Here, we adapt this approach to continuous variation in polygenic traits.

In the endemic equilibrium, the number of susceptible individuals of each type, say i , should not change over time (apart from random fluctuation). Thus, for each type i , the number of newly infected susceptibles should be equal to the number of recovering infecteds,

$$c\gamma_i\bar{\varphi}_{\text{inf}} S_i(t) \frac{I}{N} = \alpha_i I_i(t) \quad (19)$$

where $S_i(t)$ is the number of susceptible individuals of type i at time t , c the contact rate, γ_i the susceptibility of type i , $\bar{\varphi}_{\text{inf}}$ the mean infectivity of the infected individuals in the endemic equilibrium, I the total number of infected individuals in the endemic equilibrium, N total population size, α_i the recovery rate for type i , and $I_i(t)$ the number of infected individuals of type i at time t . The left-hand side in Equation (19) represents the decrease in the number of susceptibles due to transmission (infection), while the right-hand side represents the increase of the number of susceptibles due to recovery of infected individuals. Our interest is in the solution of Equation (19) for I_i (or equivalently, for $S_i = N_i - I_i$). Above, we used i to index individuals. Here, we use i also to index types, since each individual will be genetically unique for polygenic traits, so that a type corresponds to an individual. Moreover, we treat S_i and I_i as noninteger because our interest is in their expectation. Note that the mean infectivity of the infected individuals in the endemic equilibrium ($\bar{\varphi}_{\text{inf}}$) will differ from the simple population average of infectivity ($\bar{\varphi}$) when infectivity is correlated to susceptibility and/or recovery.

Equation (19) can be solved for the endemic prevalence in type i , $P_i = I_i/N_i$, N_i denoting the total number of individuals of type i in the population, and substituting $N_i = S_i + I_i$,

$$P_i = \frac{\mathcal{R}_{0,i} P}{\mathcal{R}_{0,i} P + 1} \quad (20a)$$

where $P = I/N$ denotes the overall endemic prevalence in the population (Equation 1), and

$$\mathcal{R}_{0,i} = \frac{c\gamma_i\bar{\varphi}_{\text{inf}}}{\alpha_i}. \quad (20b)$$

Equations (20a) and (20b) make no assumptions on the distribution of γ , φ , and α , and are thus not restricted to log-normal distributions. Although Equation (20b) is similar to Equation (8), note that $\mathcal{R}_{0,i}$ differs from the genotypic value for R_0 ($G_{R_0,i}$; We use a symbol slightly different from R to highlight this difference). The $\mathcal{R}_{0,i}$ is a function of the mean infectivity of the infected individuals in the endemic equilibrium ($\bar{\varphi}_{\text{inf}}$), while $G_{R_0,i}$ is a function of the infectivity of the individual itself (φ_i). Our interest here is in the prevalence for an individual with susceptibility γ_i and recovery rate α_i in the endemic equilibrium, where i is exposed to the mean infectivity of the infected individuals. For this reason, $\mathcal{R}_{0,i}$

is a function of $\bar{\varphi}_{\text{inf}}$ rather than φ_i . The $G_{R_0,i}$, in contrast, defines the contribution of an individual's genes to R_0 (Equation 18), which is relevant for response to selection. Note that $\bar{\varphi}_{\text{inf}}$ depends on the multivariate distribution of γ , φ , and α .

To find the endemic prevalence, we need to solve Equations (20a) and (20b) for P . While we found an approximate analytical solution for the case without (correlated) genetic variation in infectivity, the resulting expression is very complex (not shown). We therefore used a numerical solution, which is easily obtained (see Appendix B for methods, and Supplementary Material 1 for an R-code). We validated the numerically obtained solution using full stochastic simulation of actual endemics, following standard methods in epidemiology. Results of these simulations confirmed the numerically obtained solutions (Appendix C).

The solutions of Equations (20a) and (20b) show that variation in susceptibility and/or recovery reduces the endemic prevalence, compared to the simple prediction based on R_0 (Equation 3). Hence, with variation in susceptibility and/or recovery, prevalence is always lower than predicted by Equation (3) (Figure 3; as expected with heterogeneity; Greenhalgh *et al.* 2000). Note that genetic variation in infectivity has no effect on the prevalence (beyond its trivial effect on the mean of G_{R_0} , Equation 18), as long as infectivity is not correlated to susceptibility and/or recovery.

Moreover, it follows from Equations (20a & b) that the effects of genetic variation on the endemic prevalence are identical for susceptibility and recovery, since the $\mathcal{R}_{0,i}$ of an individual depends on the difference between its breeding values for log-susceptibility and log-recovery, $A_{\gamma,i} - A_{\alpha,i}$,

$$\mathcal{R}_{0,i} = c\bar{\varphi}_{\text{inf}} e^{A_{\gamma,i} - A_{\alpha,i}}. \quad (21)$$

Hence, with a log-normal distribution of susceptibility and recovery, and in the absence of correlated variation in infectivity, the equilibrium prevalence depends only on (R_0 and) the variance of this difference,

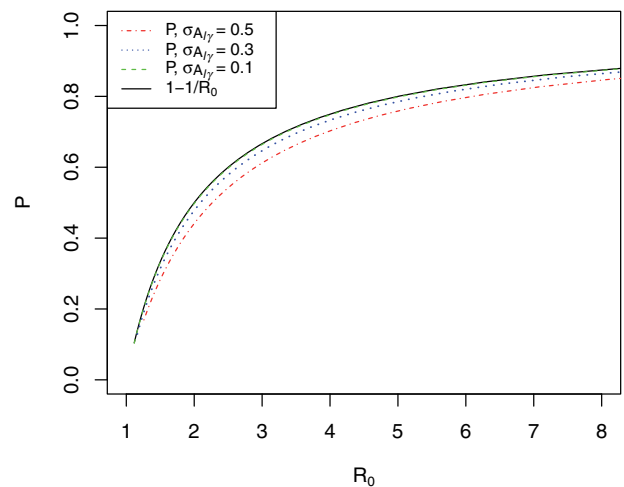


Figure 3 The impact of heterogeneity on the endemic prevalence. The solid line shows the prevalence predicted from Equation (3). The three other lines show the true prevalence (from numerically solving Equations 20a and 20b), for three levels of the additive genetic standard deviation in log susceptibility (σ_{A_γ}), and no genetic variation in infectivity or recovery. The line for $\sigma_{A_\gamma} = 0.1$ is almost identical to the solid line. Note that identical results would have been obtained with the same amount of heterogeneity in the recovery rate, or more generally in $A_{\gamma,i} - A_{\alpha,i}$, instead of in susceptibility.

$$\text{var}(A_{ly} - A_{lx}) = \sigma_{A_{ly}}^2 - 2\sigma_{A_{ly}A_{lx}} + \sigma_{A_{lx}}^2. \quad (22)$$

Figure 3 illustrates the impact of heterogeneity on the endemic prevalence for a limited number of scenarios with genetic variation in susceptibility only. For $\sigma_{A_{ly}} = 0.1$, the effect of heterogeneity is imperceptible. For $\sigma_{A_{ly}} = 0.3$, the true prevalence is up to 2 percent point lower than the value from Equation (3). For $\sigma_{A_{ly}} = 0.5$, true prevalence is up to 6 percent point lower. This maximum difference occurs at a contact rate of two. Moreover, when $c=2$ and there is no variation in infectivity, prevalence is always equal to $1 - 1/c = 0.5$, irrespective of the genetic variation in susceptibility and recovery. (This is not visible in Figure 3, because the x-axis shows R_0 rather than c). This occurs because the two opposing effects mentioned at the beginning of this paragraph exactly cancel each other. More detailed results can be found in [Supplementary Material 2](#).

Genotypic value for individual binary infection status

In the previous two sections, we have considered the population-level effects of genetic heterogeneity. In the next two sections, we move to the individual level. This section focusses on the effects of an individual's genes on its own infection status, while the next section focusses on the effect of an individual's genes on the prevalence in the population.

By definition, the genotypic value for binary disease status is the expected infection status of an individual given its genotype,

$$G_{y,i} = E(y_i | A_{ly,i}, A_{lx,i}). \quad (23)$$

Thus, the G_y represents the direct genetic effect (DGE) on the own phenotypic value (including the mean, \bar{y} , here; Note the distinction between subscript y , indicating individual binary infection status, and γ , indicating susceptibility). The genotypic value of an individual is not a function of its breeding value for log-infectivity, since an individual's infectivity does not affect its own infection status. Hence, Equation (23) does not condition on $A_{ly,i}$.

In the previous section, we used Equations (20a) and (20b) to investigate the effect of heterogeneity on the endemic prevalence in the population. Equation (20a) shows the expected prevalence of an individual of type i . However, since prevalence is simply the mean of binary infection status, Equation (20a) may also be interpreted as the expected phenotypic value for infection status ($y=0,1$) of an individual, given its genotype (specifically, the γ_i and α_i components of $\mathcal{R}_{0,i}$). Hence, Equation (20a) also represents the genotypic value for binary infection status,

$$G_{y,i} = \frac{\mathcal{R}_{0,i}P}{\mathcal{R}_{0,i}P + 1} \quad (24)$$

where $\mathcal{R}_{0,i}$ follows from Equation (20b). The same result was found by [Bijma \(2020\)](#), but based on a different approach. Thus, the G_y refers to the expected binary infection status of individual i , conditional on its genotype, in a population with prevalence P . Equations (21) and (24) imply that susceptibility and recovery are equally important for the infection status of an individual. For example, an individual with $A_{ly} = -0.1$ has the same expected infection status as an individual with $A_{lx} = +0.1$.

Calculation of G_y from Equation (24) requires knowledge of the endemic prevalence P . In the previous section, we used a numerical approach to find P , because our interest was in the effects of heterogeneity on P . In applied breeding, however, breeders may often have a reasonable idea of realistic values for the endemic

prevalence, and a numerical solution may not be needed to find G_y , or have little added value. (The dependence of the breeding value for binary infection status on the endemic prevalence will be given in Equation 34 below).

Validation

We used stochastic simulation of endemics, following standard methods in epidemiology ([Appendix C](#)), to validate Equation (24). Figure 4, A–C shows the mean observed infection status of individuals as a function of their genotypic value G_y . For all three panels in Figure 4, regression coefficients were very close to 1, showing that G_y is an unbiased linear predictor of individual infection status.

We numerically investigated the relative amount of nonadditive genetic variance in G_y by comparing the full variance in G_y (with G_y calculated from Equation 24) with the variance explained by linear regression of G_y on A_{ly} , using simulated data with variation in susceptibility only. (Note that simulation of variance in recovery would give identical results, as can be inferred from Equation 21). Results (not shown) revealed only little nonadditive genetic variance. For example, for $P=0.2$ and $\sigma_{A_{ly}}^2 = 0.5^2$, more than 96% of the genotypic variance in y was additive. Thus, the breeding value for own infection status is very similar to the genotypic value,

$$A_y \approx G_y - P, \quad (25)$$

where the “ $-P$ ” term simply reflects subtraction of the average, $\bar{G}_y = P$, so that the mean breeding value is zero by definition. We defer further investigation of the breeding value and the additive genetic variance for individual infection status to the next section, to facilitate comparison with the corresponding measures for endemic prevalence.

Individual genetic effects for the endemic prevalence

The previous section focused on the genetic effects of individuals on their own infection status. In this section, we will consider the genetic effects of individuals on the endemic prevalence in the population. In other words, the previous section focused on the contribution of genetic effects to the variation in infection status among individuals, while this section considers the genetic effects that are relevant for response to selection. We will present expressions for the genotypic value, breeding value and additive genetic variance for the endemic prevalence. The genotypic value will reflect the full genetic effect of an individual on the endemic prevalence in the population, while the breeding value reflects the additive component thereof. The last part of this section contains a comparison of the breeding value for endemic prevalence and that for individual infection status.

The relationship between R_0 and the endemic prevalence (Equation 3) suggests we can translate the individual genotypic value for R_0 (Equations 8 and 9) to the scale of prevalence, by defining an individual genotypic value for prevalence as

$$G_{P,i} = 1 - \frac{1}{G_{R_0,i}} \quad (26a)$$

Because this definition is based on Equation (3), it ignores the effect of heterogeneity on the relationship between P and R_0 (Figure 3). We will investigate the relevance of this approximation numerically in the section on response to selection. Substituting Equation (9) into Equation (26a) yields an expression for the

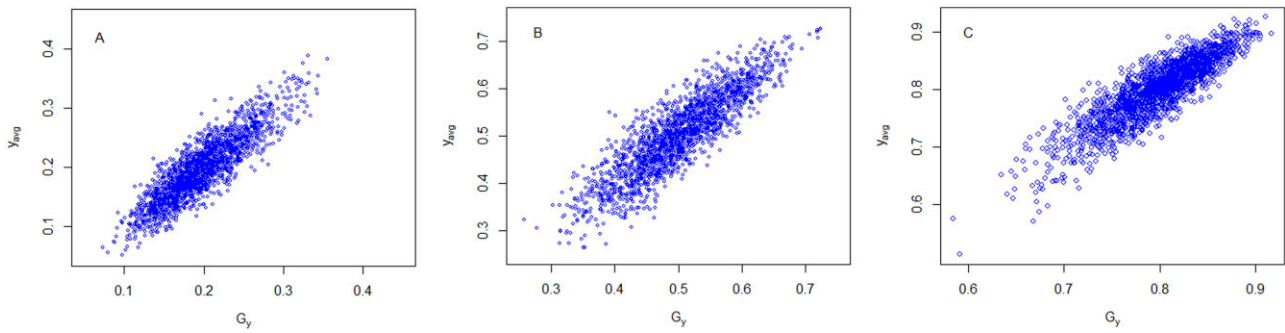


Figure 4 Validation of the genotypic value for individual binary infection status. Panels show a scatter plot of the mean observed infection status of individuals (y-axis) as a function of their genotypic value for infection status (G_y , x-axis, Equation 24). For genetic variation in susceptibility only, with $\sigma_{A_{I_0}}^2 = 0.32$ $N = 2000$ individuals, and a total of 300,000 events (sum of recoveries and infections). (A) $P = 0.2$; $b_{y,G_y} = 0.997$. (B) $P = 0.5$, $b_{y,G_y} = 1.007$. (C) $P = 0.8$, $b_{y,G_y} = 1.005$.

genotypic value of an individual for the endemic prevalence, in terms of its breeding value for the logarithm of R_0 ,

$$G_{P,i} = 1 - e^{-\ln(c) - A_{IR_0,i}} \quad (26b)$$

Because the term in the exponent is normally distributed, $1 - G_P$ follows a log-normal distribution,

$$(1 - G_P) \sim \log N\left[\mu = -\ln(c), \sigma^2 = \sigma_{A_{IR_0}}^2\right] \quad (27)$$

The mean and variance of the genotypic values for prevalence, therefore, follow from the properties of the log-normal distribution,

$$E(G_P) = 1 - c^{-1} e^{\frac{1}{2}\sigma_{A_{IR_0}}^2} \quad (28)$$

$$\text{var}(G_P) = c^{-2} \left(e^{2\sigma_{A_{IR_0}}^2} - e^{\sigma_{A_{IR_0}}^2} \right) \quad (29)$$

To enhance the interpretation of Equation (29), we can express it as a function of R_0 or of the endemic prevalence. Substituting Equation (18) into Equation (29) yields an expression for the genetic variance in prevalence as a function of R_0 ,

$$\text{var}(G_P) = \frac{1}{R_0^2} \left(e^{3\sigma_{A_{IR_0}}^2} - e^{2\sigma_{A_{IR_0}}^2} \right) \quad (30a)$$

Given the definition of G_P in Equation (26a), Equation (30a) is exact. Next, substituting Equation (3) into Equation (30a) yields an expression for $\text{var}(G_P)$ as a function of the endemic prevalence,

$$\text{var}(G_P) \approx (1 - P)^2 \left(e^{3\sigma_{A_{IR_0}}^2} - e^{2\sigma_{A_{IR_0}}^2} \right) \quad (30b)$$

Equation (30b) is approximate, because the relationship between P and R_0 given in Equation (3) is approximate with heterogeneity. Equations (30a) and (30b) show how the genotypic variance for endemic prevalence depends on the level of R_0 or equivalently, on the level of the endemic prevalence. Hence, in contrast to ordinary additive genetic traits, the genetic variance for endemic prevalence is a function of the level of the endemic prevalence (Equation 30b).

Figures 5A and B illustrate that the standard deviation in genetic values for endemic prevalence is considerably larger at lower R_0 , or equivalently, at lower prevalence. Hence, even though the genetic variance in R_0 decreases with the level of R_0 (Figure 2), the genetic variance for prevalence increases strongly when R_0 decreases. This result originates from the increasing slope of the relationship between prevalence and R_0 when R_0 decreases (Figure 1). In other words, an equal change in R_0 has much greater impact on the endemic prevalence at low R_0 than at high R_0 , which is well-known in epidemiology (e.g., Metz 1978; Bolker and Grenfell 1996). Hence, for a constant variance in the breeding value for the logarithm of R_0 , the genetic variance for endemic prevalence is much greater at lower prevalence. Moreover, genetic selection for lower prevalence will lead to an increase in the genetic variance for prevalence.

Figure 6 shows some examples of the distribution of the genotypic value for endemic prevalence, for different values of R_0 and the corresponding endemic prevalence. For the scenarios in Figure 6, the observed-scale heritability of individual infection status does not exceed 0.022 (see Figure 7 below). The panels illustrate that the genotypic standard deviation for endemic prevalence is relatively large, particularly when prevalence is small. For example, for $R_0 = 1.67$ ($P \approx 0.4$; Panel B), the standard deviation in genotypic values for prevalence is around 0.19 (see also Figure 5), and values between ~ 0 and ~ 0.7 are quite probable. Hence, despite the low observed-scale heritability of individual infection status, the probable values of G_P span as much as 70% of the full 0-1 range of endemic prevalence.

Breeding value and additive genetic variance for prevalence

The genotypic value for prevalence is not identical to the additive genetic value (i.e., breeding value) for prevalence, because the exponential function in Equation (26b) is nonlinear, so that G_P contains a nonadditive component. Appendix D shows that the linear regression coefficient of G_P on A_{IR_0} is equal to $c^{-1} \exp\left(\frac{1}{2}\sigma_{A_{IR_0}}^2\right)$. Substituting $c^{-1} = e^{\sigma_{A_{IR_0}}^2}/R_0$ (from Equation 18), shows that the breeding value for prevalence is given by

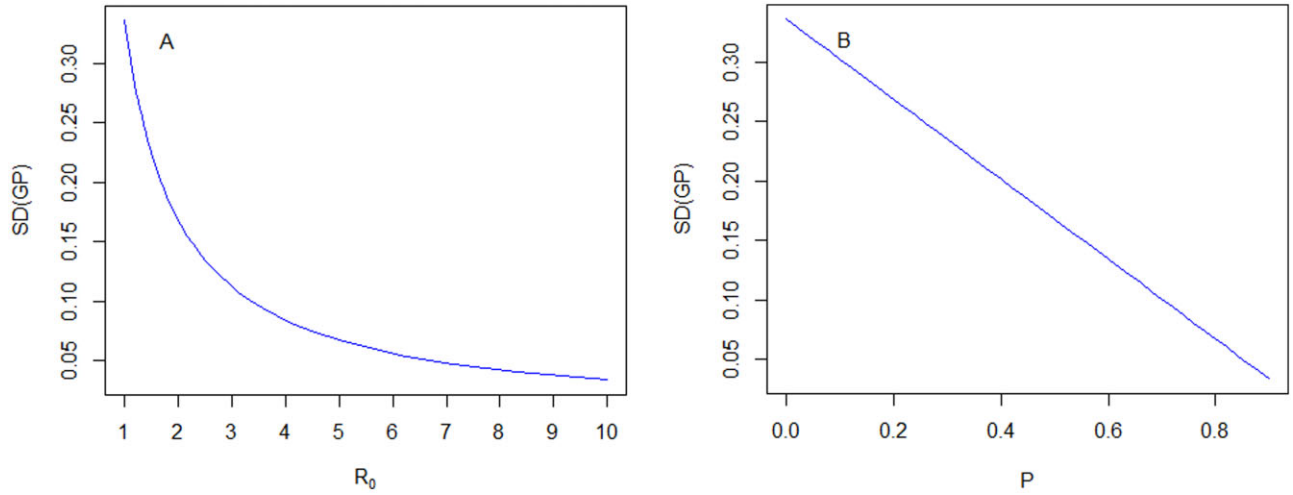


Figure 5 Genetic standard deviation for endemic prevalence as a function of R_0 (A), and as a function of the level of the endemic prevalence (B). From Equations (30a) and (30b). For $\sigma_{A|R_0}^2 = 0.3^2$. In (A), x-axis values below $R_0 = 1$ are omitted, because equilibrium prevalence is zero (the infection is not present) and Equation (30a) does not apply.

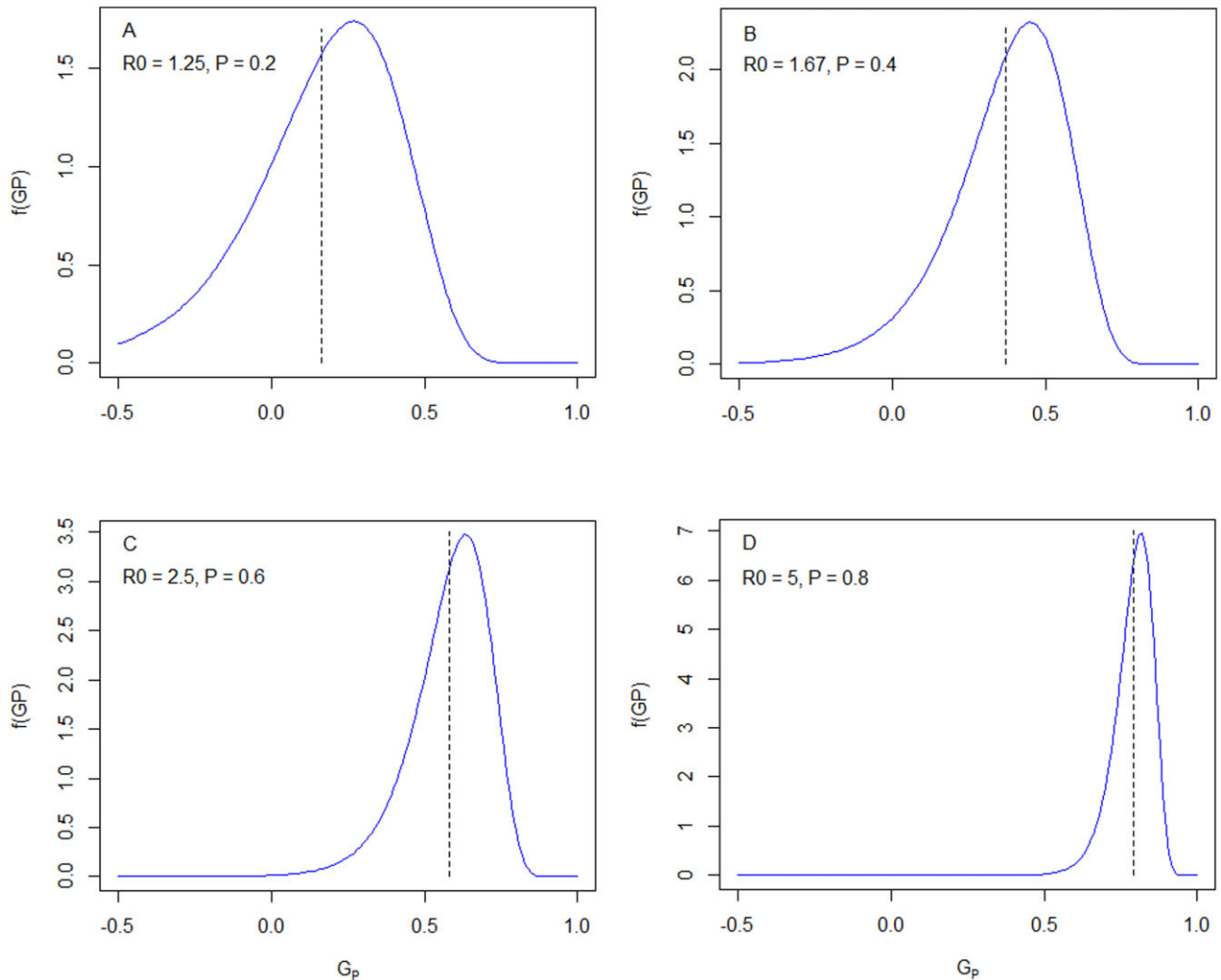


Figure 6 Distribution of individual genotypic values for prevalence (G_P), for different values of R_0 , or equivalently, different (approximate) values of the endemic prevalence. (A) $R_0 = 1.25, P = 0.2$. (B) $R_0 = 1.67, P = 0.4$. (C) $R_0 = 2.5, P = 0.6$. (D) $R_0 = 5, P = 0.8$. The distribution is given by

$$f(G_P) = \frac{1}{(1-G_P)\sigma_{A|R_0}\sqrt{2\pi}} \exp\left(-\frac{(\log(1-G_P)+\log(c))^2}{2\sigma_{A|R_0}^2}\right),$$
 with domain $G_P = (-\infty, 1)$. For $\sigma_{A|R_0}^2 = 0.3^2$. The dashed vertical line shows the mean of G_P . Note that G_P can take negative values while prevalence cannot. This is because G_P reflects the genetic effect of an individual on the prevalence of the population, not the expected value of its own infection status. Thus, negative values for G_P are possible, as long as P is positive. Note that P is very close to the average of the distributions shown.

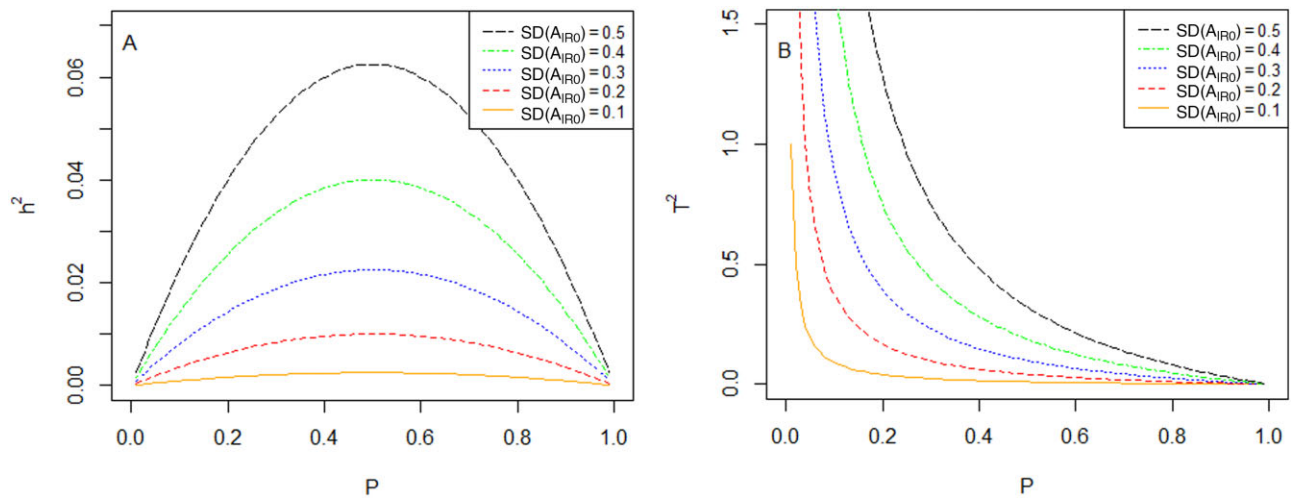


Figure 7 (A) Observed-scale heritability (h^2_y) of individual binary infection status ($y = 0/1$) as a function of the endemic prevalence, for different additive genetic standard deviations in the logarithm of R_0 ($SD(A_{IR0})$) (from Equation 36). (B) Ratio of additive genetic variance for prevalence and phenotypic variance in infection status (T^2_p), as a function of the endemic prevalence. From Equations (37) and (32a). In both panels, there is no genetic variation in infectivity.

$$A_{P,i} = \frac{1}{R_0} e^{\sigma_{A_{IR0}}^2} A_{IR0,i} \quad (31a)$$

Given the definition of G_P in Equation (26a), this result is exact. With limited heterogeneity, this result is approximately equal to

$$A_{P,i} \approx \frac{1}{R_0} A_{IR0,i} \quad (31b)$$

Thus, the breeding value for prevalence is proportional to the reciprocal of R_0 . The additive genetic variance for prevalence equals

$$\sigma_{A_P}^2 = \frac{1}{R_0^2} e^{2\sigma_{A_{IR0}}^2} \sigma_{A_{IR0}}^2 \quad (32a)$$

and, with limited heterogeneity,

$$\sigma_{A_P}^2 \approx \frac{1}{R_0^2} \sigma_{A_{IR0}}^2 \quad (32b)$$

or, expressed as a function of endemic prevalence,

$$\sigma_{A_P}^2 \approx (1 - P)^2 \sigma_{A_{IR0}}^2 \quad (32c)$$

Equation (32c) shows that the additive genetic variance in endemic prevalence increases strongly when prevalence decreases, similar to the relationship between the genotypic variance and endemic prevalence (Figure 5). This result suggests that response of endemic prevalence to selection will be greater at lower levels of the prevalence, which we will further investigate below in the section on response to selection.

The relative amount of nonadditive genetic variance in the endemic prevalence is determined by the magnitude of $\sigma_{A_{IR0}}^2$ (Appendix D). For realistic values of $\sigma_{A_{IR0}}^2$, the vast majority of the genotypic variance in prevalence is additive. For example, for $\sigma_{A_{IR0}}^2 = 0.5^2$, 88% of the variance in G_P is additive. Hence, the

distinction between the breeding value for prevalence (A_P) and the genotypic value for prevalence (G_P) seems of minor importance, and results in Figures 5 and 6 will closely resemble those for the additive genetic effects.

Breeding value and heritability for infection status vs prevalence

Appendix E shows that, in the absence of genetic variation in infectivity, the breeding value for endemic prevalence is approximately a factor $1/P$ greater than the breeding value for individual infection status,

$$A_{P,i} \approx \frac{1}{P} A_{y,i} \quad (33)$$

Note, in contrast to genotypic values, breeding values are expressed as a deviation from their mean here. The A_y is the ordinary observed-scale breeding value for binary infection status that breeders are familiar with.

Equation (33) implies that the impact of an individual's genes on the response of the endemic prevalence to selection is considerably larger than their impact on the infection status of the individual itself, particularly when the endemic prevalence is small. Consider, for example, an individual with $A_{y,i} = -0.02$ in a population with an endemic prevalence of $P = 20\%$. The expected infection status of this individual in the current population equals $0.20 - 0.02 = 0.18$. Hence, on average, this individual will be infected 18% of the time. However, its breeding value for prevalence equals $A_{P,i} = -0.02/0.2 = -0.10$. Hence, if we select individuals with $A_{y,i} = -0.02$ as parents of the next generation, then the endemic prevalence will go down to $0.20 - 0.10 = 0.10$. In other words, the response to selection will be fivefold greater than suggested by the ordinary breeding value for individual infection status (since $1/P = 1/0.2 = 5$). We will numerically validate this theoretical result in the section on response to selection below.

The relationship between the breeding value for prevalence and the breeding value for own infection status shown in

Equation (33) suggests a relatively simple expressions for A_y . Such a simple expression would be convenient, because the alternative is to calculate A_y from Equation (25), which requires solving Equations (20a) and (20b) numerically. On combining Equations (31a) and (33), and assuming limited heterogeneity, so that $e^{\sigma_{A_{IR_0}}} \approx 1$ and $1/R_0 \approx (1 - P)$, the breeding value for individual infection status becomes

$$A_{y,i} \approx P(1 - P)A_{IR_0,i} \quad (34)$$

We used stochastic simulation to validate this expression and investigate its precision. Results show that Equation (34) closely matches the regression of individual binary infection status on the breeding value for the logarithm of R_0 for realistic levels of heterogeneity ($\sigma_{A_{IR_0}}^2 \leq 0.5^2$; Supplementary Material 3; the good fit results from compensating errors due to the approximations). Hence, Equation (34) is sufficiently precise for practical purposes. Note that, since infectivity does not affect the infection status of an individual itself, a potential component due to infectivity has to be left out of the $A_{IR_0,i}$ term when calculating Equation (34). In other words, in Equation (34) the $A_{IR_0,i}$ should include only the breeding values for the logarithm of susceptibility and recovery (see Equation 10a).

It follows from Equation (34) that the additive genetic variance in individual binary infection status equals

$$\sigma_{A_y}^2 \approx P^2(1 - P)^2 \sigma_{A_{IR_0}}^2. \quad (35)$$

Next, the observed-scale heritability of binary infection status follows from dividing Equation (35) by the phenotypic variance of binary infection status, $\sigma_y^2 = P(1 - P)$, giving

$$h_y^2 \approx P(1 - P) \sigma_{A_{IR_0}}^2. \quad (36)$$

Hence, the observed-scale heritability for binary infection status has a maximum at a prevalence of 0.5, and goes to zero at a prevalence of zero or one, just like the heritability of binary phenotypes for noncommunicable polygenic traits (Robertson 1950; Figure 7A; assuming the infinitesimal model at the level of the logarithm of R_0 , so that $\sigma_{A_{IR_0}}^2$ is constant).

The ratio of additive genetic variance for prevalence over phenotypic variance in binary infection status is given by,

$$T_P^2 = \frac{\sigma_{A_P}^2}{P(1 - P)} \approx \frac{1 - P}{P} \sigma_{A_{IR_0}}^2 \quad (37)$$

with $\sigma_{A_P}^2$ taken from Equation (32a) or Equation (32c). The T_P^2 is an analogy of heritability, but the numerator represents the additive genetic variance relevant for response to selection in endemic prevalence, rather than for individual binary infection status. The T_P^2 , therefore, reflects the genetic variance that can be used for response to selection, whereas h_y^2 reflects the relative contribution of additive genetic effects to the phenotypic variance in binary infection status (Bijma et al. 2007; Bijma 2011).

Figures 7A and B show a comparison of h_y^2 and T_P^2 for a population without genetic variation in infectivity, with genetic variances in the logarithm of R_0 ranging from 0.1^2 through 0.5^2 . In Figure 7A, the maximum value of h_y^2 equals 0.0625, for $P = 0.5$ and $\sigma_{A_{IR_0}}^2 = 0.5^2$. Given that genetic variances greater than $\sigma_{A_{IR_0}}^2 = 0.5^2$ are very large (as argued above), observed-scale heritabilities of binary infection status greater than ~ 0.06 are unlikely for endemic infectious diseases. The heritabilities in Figure 7A agree

with the findings of Hulst et al. (2021), who used stochastic simulation of actual endemics and analysis of the resulting binary infection status data with a linear animal model. Figure 7B shows that T_P^2 increases strongly when prevalence goes down. Figure 7 illustrates that T_P^2 and h_y^2 differ by a factor of approximately P^2 , so that the additive genetic variance in prevalence is (much) greater than the additive genetic variance in individual infection status, and may even exceed the phenotypic variance at low values of the endemic prevalence (i.e., $T_P^2 > 1$).

In conclusion, in this section, we have presented expressions for the breeding value for prevalence (Equation 31) and for individual infection status (Equation 34), and for the corresponding genetic variances. With realistic levels of heterogeneity, the breeding value for prevalence is a factor $1/P$ greater than the breeding value for individual infection status. This result suggests that response to selection should be considerably greater than expected based on ordinary heritability of individual infection status. We will test this hypothesis in the next section.

Response to selection

The higher genetic variance for prevalence at lower values of the prevalence (Equations 30 and 32, Figures 5B and 6) suggests that the response of the endemic prevalence to selection should increase when the prevalence decreases. To validate and illustrate this hypothesis, we stochastically simulated an endemic infectious disease in a large population undergoing mass selection for individual infection status. Hence, the individuals with the lowest observed average infection status were selected as parents of the next generation. Simulations were based on standard methods in epidemiology, not making use of the above theory (Appendix F).

Figure 8A shows the observed prevalence (i.e., the mean binary infection status in each generation), the mean breeding value for prevalence and the mean breeding value for binary infection status, for ~ 70 generations of selection. Response in prevalence increases strongly when prevalence decreases, and the infection disappears in the final generation. There is excellent agreement between the observed prevalence and the breeding value for prevalence, showing that the change in \bar{A}_P indeed predicts the change in prevalence. In contrast, the response in prevalence deviates substantially from the response in the breeding value for individual infection status (\bar{A}_y), particularly at lower values of the prevalence. Hence, while the breeding value for infection status correctly predicts the average individual infection status within a generation (Figure 4), the change in \bar{A}_y considerably underestimates the response to selection. Furthermore, given the weak selection and the low value of the observed-scale heritability of binary infection status, which did not exceed 0.022 in Figure 8A, response to selection in prevalence is remarkably large, unless prevalence is high. This result agrees with findings of Hulst et al. (2021).

We also investigated the prediction of response to mass selection with a very simple expression assuming linearity, with genetic variation in susceptibility only, and only for $\sigma_{A_{Iy}}^2 = 0.3^2$. The response of prevalence to selection follows from the ordinary breeder's equation, applied to prevalence,

$$R_P = i \rho_{A_P, SC} \sigma_{A_P} \quad (38a)$$

where i is the intensity of selection, $\rho_{A_P, SC}$ the accuracy of selection, being the correlation between the selection criterion and the true breeding value for prevalence in the candidates for selection, and σ_{A_P} the additive genetic standard deviation for prevalence.

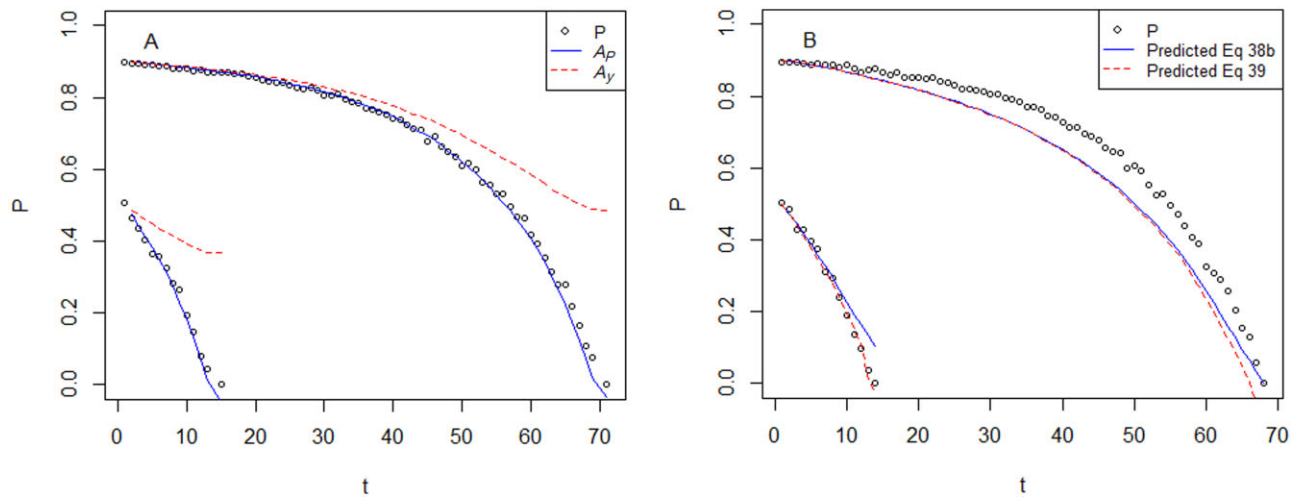


Figure 8 Response to selection in prevalence for 70 generations of mass selection of the host population. (A) compares observed prevalence to observed true breeding values for prevalence and for individual disease status. (B) compares observed prevalence to predicted prevalence. For two populations, one starting at a prevalence of ~90% ($c = 10$), the other starting at a prevalence of ~50% ($c = 2$). Each generation, the 50% individuals with the lowest average infection status were selected as parents of the next generation. With genetic variation in susceptibility only, and $\sigma_{A_{Y_i}}^2 = 0.3^2$. For a population of $N = 4,000$ individuals, a total of 15,000 events (sum of infections and recoveries) per generation, consisting of a burn-in of 10,000 events and 5,000 recorded events. Hence, selection is based on 1.25 events per individual on average, indicating a limited amount of phenotypic data. Observed-scale heritability for binary infection status in any generation can be read from Figure 7A using an x-axis value corresponding to the prevalence in that generation. (A) Observed prevalence (circles), breeding value for prevalence (\bar{A}_P , solid blue line) and breeding value for individual infection status (\bar{A}_Y , dashed red line). Results for breeding values are the cumulative change in breeding value in each generation plus the initial prevalence. Breeding values were taken from Equation (31a) for \bar{A}_P and from Equation (34) for \bar{A}_Y . (B) Predicted (lines) versus observed (circles) prevalence. Prevalence was predicted from Equation (38b) (blue solid line) or Equation (39) (red dashed line).

However, at present breeders are not familiar with genetic parameters for prevalence. For this reason, we based our predictions on genetic parameters for binary infection status, because these are typically available to breeders. Because the breeding value for prevalence and the breeding value for individual infection status differ by a factor $1/P$ when there is no genetic variation in infectivity (Equation 33), we simply upscaled the response to mass selection predicted for individual infection status from the breeder's equation (Walsh and Lynch 2018) by this factor, giving

$$R_P = \iota \rho_{A_Y, \bar{Y}} \sigma_{A_Y} \frac{1}{P} \quad (38b)$$

where the intensity of selection, defined as the standardized selection differential in mean individual infection status, was $\iota = (\bar{Y}_{\text{selected}} - \bar{Y})/\sigma_{\bar{Y}}$, the accuracy $\rho_{A_Y, \bar{Y}}$ is the correlation between the selection criterion (\bar{Y}_i here) and the true breeding value for individual infection status, and P is the prevalence in the generation of the selection candidates. Thus, the selection criterion was the observed average infection status of the individual itself (\bar{Y}_i , mass selection). Hence, the numerator of Equation (38b) represents the predicted response for individual binary infection status, which is multiplied by a factor $1/P$ to find response in prevalence. To implement Equation (38b), we calculated the $\rho_{A_Y, \bar{Y}}$ as the observed correlation between the true breeding values for binary infection status (A_Y ; Equation 34) and the selection criterion (\bar{Y}) in the candidates for selection. Hence, we did not attempt to predict the accuracy of selection. (Note that Equation 38b also applies to selection on EBV for infection status, where the accuracy refers to the accuracy of these EBV, being the correlation between the true and the estimated breeding value for infection status, ρ_{A_Y, \hat{A}_Y}).

Figure 8B shows a comparison of observed and predicted prevalence. Above a prevalence of ~0.5, response predicted from Equation (38b) is somewhat larger than observed response, while the reverse is true below a prevalence of ~0.5 (Note, response to selection in a generation is reflected by the slope of the figure). Nevertheless, agreement between observed and predicted response is remarkably good given the very unrealistic assumption of linearity in Equation (38b) (i.e., bivariate normality of A_Y and \bar{Y}). Because selection was based on mean individual infection status recorded over a period lasting on average only 1.25 events per individual (see legend Figure 8), many values were either 0 or 1, implying strong deviations from normality.

When the prevalence was smaller than 0.5, response to selection was quite large. Hence, there was a meaningful difference in prevalence between parent and offspring generations. Because the P in Equation (38b) refers to the prevalence in the parent generation, while response is realized in the offspring generation, Equation (38b) resulted in underprediction of response to selection when response was large. This underprediction disappeared when using prevalence in the offspring generation in the $1/P$ term in Equation (38b). However, because prevalence in the offspring generation is initially unknown, as it depends on the response to selection, this prediction required solving the expression $R = \iota \rho_{A_Y, \bar{Y}} \sigma_{A_Y} / (P + R)$, yielding

$$R_P = \frac{1}{2} \left(-P + \sqrt{P^2 + 4\iota \rho_{A_Y, \bar{Y}} \sigma_{A_Y}} \right). \quad (39)$$

For a prevalence smaller than ~0.5, predictions from Equation (39) were very close to the observed response in prevalence [Figure 8B; for $P > 0.5$, results of Equations (38b) and (39) are almost identical].

In conclusion, results in this section show that response to selection in the prevalence of endemic infectious diseases is a factor $1/P$ greater than suggested by the ordinary breeding values for individual binary infection status. Thus, breeders can predict response to selection by upscaling the selection differential in the usual estimated breeding values for binary infection status by a factor $1/P$.

Direct and indirect genetic variance for endemic prevalence

In this section, we partition the total additive genetic variance for endemic prevalence into direct and indirect genetic components. This partitioning is relevant, because IGE respond fundamentally different to selection than DGE (Griffing 1967, 1977; Wright 1985; Moore et al. 1997; Muir 2005; Bijma 2010, 2011; see Discussion). We can partition the total breeding value for prevalence into a direct and an indirect component,

$$A_P = A_{P_D} + A_{P_I} \quad (40)$$

Analogously, we can partition the full additive genetic variance in prevalence into components due to direct genetic variance, indirect genetic variance and a covariance,

$$\sigma_{A_P}^2 = \sigma_{A_{P_D}}^2 + 2\sigma_{A_{P_D}A_{P_I}} + \sigma_{A_{P_I}}^2 \quad (41)$$

In the absence of genetic variation in infectivity, the breeding value for own infection status is a fraction P of the breeding value for prevalence (Equation 33). Hence, a fraction P of the additive genetic effects of susceptibility and recovery on prevalence affects the infection status of the individual itself and thus represents a direct effect, while the remaining fraction $(1 - P)$ represents an indirect effect. For infectivity, the entire genetic effect is indirect, because an individual's infectivity does not affect its own infection status. It follows from Equations (31a) that

$$A_{P_D} = \frac{e^{\sigma_{A_{IR_0}}^2}}{R_0} P (A_{I_Y} - A_{I_X}) \quad (42a)$$

$$A_{P_I} = \frac{e^{\sigma_{A_{IR_0}}^2}}{R_0} ((1 - P)A_{I_Y} + A_{I_\phi} - (1 - P)A_{I_X}) \quad (42b)$$

Note that Equation (42a) represents the breeding value for individual infection status (A_Y), but the current expression emphasizes the partitioning of A_P into direct and indirect effects. The direct and indirect genetic (co)variances are given by

$$\sigma_{A_{P_D}}^2 = \frac{e^{2\sigma_{A_{IR_0}}^2}}{R_0^2} P^2 (\sigma_{A_{I_Y}}^2 - 2\sigma_{A_{I_Y}A_{I_X}} + \sigma_{A_{I_X}}^2) \quad (43a)$$

$$\sigma_{A_{P_D}A_{P_I}} = \frac{e^{2\sigma_{A_{IR_0}}^2}}{R_0^2} P \left\{ (1 - P) (\sigma_{A_{I_Y}}^2 - 2\sigma_{A_{I_Y}A_{I_X}} + \sigma_{A_{I_X}}^2) + \sigma_{A_{I_Y}A_{I_\phi}} - \sigma_{A_{I_Y}A_{I_X}} \right\} \quad (43b)$$

$$\begin{aligned} \sigma_{A_{P_I}}^2 = & \frac{e^{2\sigma_{A_{IR_0}}^2}}{R_0^2} \left\{ (1 - P)^2 (\sigma_{A_{I_Y}}^2 - 2\sigma_{A_{I_Y}A_{I_X}} + \sigma_{A_{I_X}}^2) \right. \\ & \left. + 2(1 - P)(\sigma_{A_{I_Y}A_{I_\phi}} - \sigma_{A_{I_\phi}A_{I_X}}) + \sigma_{A_{I_\phi}}^2 \right\} \end{aligned} \quad (43c)$$

Figure 9 shows the total additive genetic variance for the endemic prevalence and the fractions due to DGE, IGE and their covariance, for a scenario with equal genetic variances in susceptibility, infectivity and recovery and covariances equal to zero. For an endemic prevalence smaller than 0.5, IGE contribute

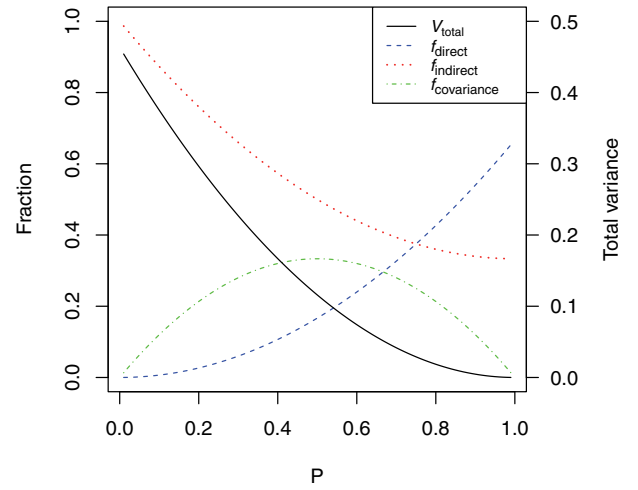


Figure 9 Total additive genetic variance in endemic prevalence (V_{total} ; secondary y-axis) and the relative contributions of DGE, IGE and their covariance (f_{direct} , f_{indirect} , and $f_{\text{covariance}}$; primary y-axis). For constant values of $\sigma_{A_{I_Y}}^2 = \sigma_{A_{I_\phi}}^2 = \sigma_{A_{I_X}}^2 = 0.3^2$ and covariances equal to zero. Results are obtained from Equations (32b), (3), and (43a)–(43c).

the majority of the genetic variance. For example, for an endemic prevalence of 0.3, the total additive genetic variance consists of 6% direct genetic variance, 66% indirect genetic variance and 28% direct-indirect genetic covariance. These results imply that IGE dominate the heritable variation and response to selection for the endemic prevalence of infectious diseases, unless prevalence is high.

Discussion

We have presented a quantitative genetic theory for endemic infectious diseases, with a focus on the genetic factors that determine the endemic prevalence. We defined an additive model for the logarithm of individual susceptibility, infectivity and rate of recovery, which results in normally distributed breeding values for the logarithm of R_0 . Next, we investigated the impact of genetic heterogeneity on the population level, for both R_0 and the endemic prevalence. Results show that, despite heterogeneity, R_0 remains equal to the mean individual genotypic value for R_0 . Subsequently, we considered genetic effects of individuals on their own infection status and on the endemic prevalence in the population. Building on the breeding value for the logarithm of R_0 , we showed that genotypic values and genetic parameters for the prevalence follow from the known properties of the log-normal distribution. In the absence of genetic variation in infectivity, genetic effects for the endemic prevalence are a factor $1/\text{prevalence}$ greater than the ordinary genetic values for individual binary infection status. Hence, even though prevalence is the simple average of individual binary infection status, breeding values for prevalence show much more variation than those for individual infection status. These results imply that the genetic variance that determines the potential response of the endemic prevalence to selection is largely due to IGE, and thus hidden to classical genetic analysis and selection. For susceptibility and recovery, a fraction $1-P$ of the full genetic effect on endemic prevalence is due to IGE, whereas the effect of infectivity is entirely due to IGE. Hence, the genetic variance that determines the potential response of the endemic prevalence to selection must be much greater than expected based on classical quantitative genetic theory, particularly at low levels of the prevalence (Figure 7). We

evaluated this implication using stochastic simulation of endemics following standard methods in epidemiology, where parents of the next generation were selected based on their own infection status (mass selection). The results of these simulations show that response to selection in the observed prevalence and in the breeding value for prevalence increases strongly when prevalence decreases, and closely matches our predictions, which supports the theoretical findings presented here.

Model assumptions

Following [Anacleto et al. \(2015, 2019\)](#), [Biemans et al. \(2019\)](#), and [Pooley et al. \(2020\)](#), we assumed a linear additive model with normally distributed effects for the logarithm of susceptibility, infectivity and recovery, leading to a normal distribution of the additive genetic values for the logarithm of R_0 (Equation 10). For complex traits, it is common to assume normally distributed genetic effects, based on the central limit theorem ([Fisher 1919](#)). Because susceptibility, infectivity and recovery act multiplicatively in the expression for R_0 , and because R_0 is nonnegative, we specified a normal distribution for its logarithm. This resulted in an additive model on the log scale, which agrees with the infinitesimal model, and also translates the $[0, \infty)$ domain of R_0 to the $(-\infty, +\infty)$ domain of the normal distribution. Hence, we assumed constant genetic parameters for the logarithm of R_0 . The same approach has been used to model genetic variation in the residual variance, which is also restricted to nonnegative values ([SanCristobal-Gaudy et al. 1998](#); [Hill and Mulder 2010](#)). The log-normal distribution of genotypic values for R_0 results in a decrease of the genetic standard deviation in R_0 with decreasing R_0 (Figure 2), which seems reasonable given the presence of a lower bound for R_0 . Moreover, the log-normal distribution for R_0 is convenient, because it results in simple expressions for the breeding value and the genetic variance for prevalence (Equations 31 and 32).

The assumption of a normal distribution for the logarithm of genotypic values for R_0 also agrees with the standard implementation of generalized linear (mixed) models (GLMM; [Nelder and Wedderburn 1972](#)). R_0 refers to an expected number of infected individuals; In other words, R_0 is the expected value of count data. In GLMM, the default link function for count data is the log-link ([McCullagh and Nelder 2019](#)). Hence, our linear model for the logarithm of R_0 also agrees with common statistical practise.

The strong increase of the genetic variance in prevalence with decreasing R_0 (Figure 5A) is not due to the assumption of lognormality of R_0 . On the contrary, the log-normal distribution results in a decrease of the genetic standard deviation in R_0 with decreasing R_0 (Figure 2). The strong increase in the genetic variance in prevalence with decreasing R_0 results from the relationship between R_0 and the prevalence in the endemic steady state (Figure 1; Equation 3), which becomes steeper when R_0 is closer to one. This relationship is very well established in epidemiology since [Weiss and Dishon \(1971\)](#); e.g., [Keeling and Rohani 2011](#).

While we defined an additive genetic model for the logarithm of R_0 , we can also find the additive genetic effect (breeding value) for R_0 itself. Using results of [Appendix D](#), the breeding value for R_0 follows from

$$A_{R_0,i} = c e^{\frac{1}{2}\sigma^2_{A_{R_0}}} A_{IR_0,i}$$

Hence, the A_{R_0} represents the additive component of the genotypic value for R_0 . However, because our model is additive on the log-scale, while the genotypic value for R_0 includes nonadditive

genetic effects, we decided to build our theory on the breeding value for the logarithm of R_0 .

Throughout this manuscript, we have assumed that the means of the breeding values on the log scale are equal to zero (Equations 7 and 10). With this assumption, the expected value of a log-normal variate simplifies to $e^{\frac{1}{2}\sigma^2}$, which we have used at many places in this manuscript. This assumption can be satisfied easily by moving the mean breeding value on the log scale into the contact rate, whenever it is not zero. Thus, this assumption does not put any restriction on the validity or generality of our work, but in the application of the results, the mean of the breeding values on the log scale should be moved to the contact rate. Suppose the mean breeding values on the log scale are given by $\mu_{A_{I_1}}$, $\mu_{A_{I_0}}$ and $\mu_{A_{I_2}}$. Then the genotypic value for R_0 is given by

$$G_{R_0,i} = c_{\mu_A \neq 0} e^{(\mu_{A_{I_1}} + A_{I_1,i}) + (\mu_{A_{I_0}} + A_{I_0,i}) - (\mu_{A_{I_2}} + A_{I_2,i})}$$

where $c_{\mu_A \neq 0}$ is the contact rate for the model parameterization where the mean breeding value on the log scale is nonzero, and A denotes a breeding value on the log scale expressed as a deviation from its mean (thus A has mean zero by definition). This model can be reparametrized into Equation (9)

$$G_{R_0,i} = c e^{A_{I_1,i} + A_{I_0,i} - A_{I_2,i}}$$

using

$$c = c_{\mu_A \neq 0} e^{\mu_{A_{I_1}} + \mu_{A_{I_0}} - \mu_{A_{I_2}}}$$

Hence, to move the mean breeding value on the log scale into the contact rate, we have to multiply the original contact rate by $\exp(\mu_{A_{I_1}} + \mu_{A_{I_0}} - \mu_{A_{I_2}})$. While both parameterizations are obviously equivalent, the second results in simpler expressions and has been used throughout this manuscript. Thus, c rather than $c_{\mu_A \neq 0}$ must be used when applying our results. This is essential, because breeding values and genetic variances for the endemic prevalence depend on the contact rate (e.g., Equations 31 and 32).

In our results for response to selection (Figure 8), we have assumed that the population has reached the endemic steady state at any time. In other words, we assumed that, after a selection, the population has reached the new endemic prevalence before the next selection takes place. Whether this assumption holds true will depend on the rate of convergence of the feedback process in the disease transmission (discussed below and illustrated in Figure 10) versus the rate of genetic improvement. [Hulst et al. \(2021\)](#) show examples of convergence to the new equilibrium. If the genetic improvement goes gradually, for example when replacing part of the animals each year like in dairy cattle, and when the pathogen survives only briefly in the environment, then the prevalence of the local population will track the gradual genetic changes in the population and the improvements predicted will be observed immediately. On the other hand, if the genetic changes are large and abrupt, like when restocking broilers or fattening pigs with a new genetic stock, and when the new stock is exposed to the infectious material from the previous stock, either because the pathogen survives in the environment or from neighboring stables or pens, then it may take some time before the full effect of genetic improvement materializes.

Nevertheless, the full effect of genetic improvement will materialize over time, also when the next selection takes place before the population has reached the new endemic prevalence due to the previous selection. Thus, incomplete convergence to the new

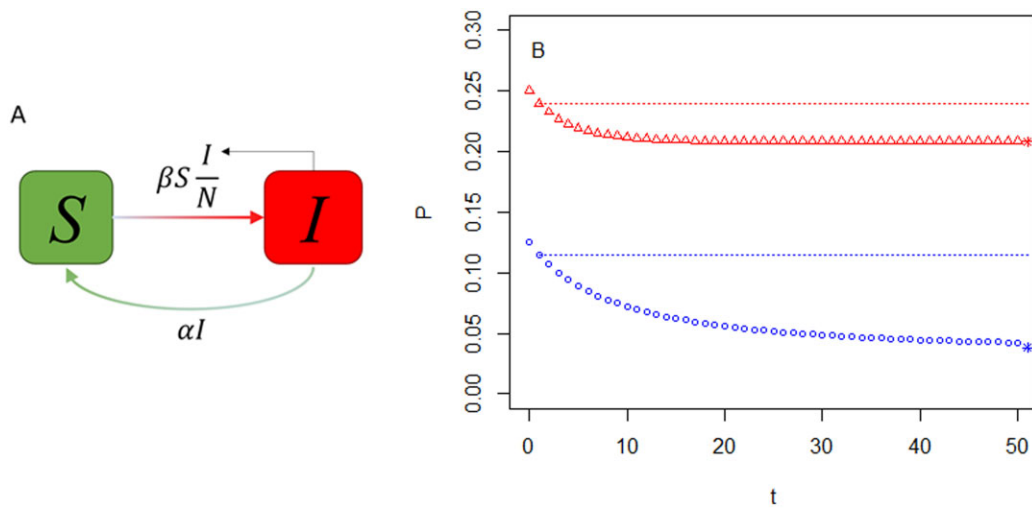


Figure 10 Positive feedback after selection for lower susceptibility. (A) Diagram of the SIS compartmental model illustrating the feedback, with the number of susceptible (S) and infectious (I) individuals and the transmission and recovery rates (ignoring heterogeneity for simplicity). A reduction in the transmission rate parameter β reduces I , which in turn reduces the transmission rate, leading to a further reduction in I , etc. (B) Convergence of the prevalence to the new equilibrium after selection. For two populations; one starting at $P = 0.25$ (red triangles; $c = 1.333$), the other at $P = 0.125$ (blue circles; $c = 1.143$). The x-axis represents cycles of the transmission loop. The horizontal dotted lines show the prevalence predicted by the breeding value for binary infection status, and represent the direct effect. The Asterisk shows the equilibrium prevalence after convergence, which occurs around $t = 22$ for the upper line, a little later than $t = 50$ for the lower line. The genetic selection differential for binary infection status equals $\Delta \bar{A}_y = -0.01$ for each of the two populations. The initial response to selection (the y-axis difference between $t = 0$ and $t = 1$) is equal to the $\Delta \bar{A}_y$ of -0.01 for both scenarios. Total response is -0.04 for the scenario with $P = 0.25$, and -0.08 for the scenario with $P = 0.125$, corresponding to $-0.01/0.25$ and $-0.01/0.125$. Results in (B) follow from iterating on Equation (20a), using a single value for $\mathcal{R}_{0,i}$, with $\alpha_i = \bar{\varphi}_{\text{inf}} = 1$ and choosing γ so that the selection differential $\Delta \bar{A}_y = -0.01$ (using $\gamma = \frac{G_y}{P(1-G_y)} / c$ from Equation 20a). In each iteration, the P in the righthand side of Equation (20a) is replaced by the P_i calculated from Equation (20a) in the previous iteration. This iteration converges to the prevalence given by Equation (3) (assuming negligible heterogeneity).

equilibrium before the next selection takes place is a transient phenomenon. It does not affect the ultimate genetic improvement, because the ultimate endemic prevalence is determined by the genetic value for R_0 , not by the previous prevalence. Incomplete convergence to the new equilibrium may actually lead to a slightly greater response, because the accuracy of selection will be a bit higher at higher prevalence, leading to a larger genetic change. In other words, the $\rho_{A_y, \bar{y}}$ in Equation (38b) will be higher, due to higher heritability (This follows from Figure 7A, where heritability is higher at higher prevalence, as long as $P < 0.5$).

Other compartmental models

In this study, we focused on endemic infectious diseases following a SIS-model, where individuals can be either susceptible (S , i.e., noninfected) or infected (I). Hence, we assumed the infection does not confer any long-lasting immunity, and we ignored the potential existence of infected classes (“compartments”) other than S and I , such as recovered infected individuals that are not yet susceptible again. Moreover, we ignored the influx of new individuals into the population due to births, and the removal of individuals due to deaths.

A key condition for validity of our results is that the pathogen can replicate only in the host individual, meaning that a reduction in infected individuals fully translates into reduced exposure of the host population to the pathogen. The mere survival of the pathogen in the environment does not violate our assumptions (see Hulst et al. 2021 for a discussion). Our conclusions are not limited to SIS models if this condition is met, but apply to all models with no longer lasting immunity. For models with temporary immunity (e.g., SIRS) or lifelong immunity (e.g., SIR) the conclusions with respect to infectivity and susceptibility will be true,

but the genetic variation in recovery may a different more restricted role.

Also, infections that do confer long-lasting immunity may show endemic behavior when a population is large enough. Measles in the human population before the introduction of vaccination are a well-known example. For such infections, the same mechanisms as discussed above will play a role and the endemic prevalence for a homogeneous population still follows from Equation (3). However, the introduction of new susceptibles by birth can no longer be ignored, and recovery of infected individuals does not result in new susceptible individuals. Thus, the role of recovery will change, and the genetic make-up of the newborn individuals becomes relevant, particularly in populations undergoing selection.

Positive feedback

The increasing difference between the breeding value for prevalence and the breeding value for individual infection status at lower prevalence (Equation 33) is a result of the increasing slope of the relationship between R_0 and the endemic prevalence (Equation 3, Figure 1). Equation (3) follows directly from a simple equilibrium condition (see text above Equation 3). However, the focus on the equilibrium partly obscures the underlying mechanism.

For a disease caused by exposure, the effect of genetic selection depends on future exposure. For an infectious disease, future exposure depends on the future number of infected individuals in the local population and on their (lifetime) infectivity, both of which are affected by the genetic selection. Thus, for infectious diseases future exposure depends on selection, leading to feedback effects. Figure 10, A and B illustrates that the difference between A_P and A_y originates from positive feedback effects in the transmission dynamics. (Figure 10 shows results for selection

against susceptibility, selection for faster recovery would yield identical results). With lower susceptibility fewer individuals will become infected, which subsequently translates into a reduced transmission rate, followed by a further reduction in the number of infected individuals, etc, resulting in a positive feedback loop (Figure 10A). The initial change in prevalence before feedback effects manifest is equal to the selection differential in breeding value for individual infection status ($\Delta\bar{A}_i$; horizontal lines in Figure 9). This change represents the direct response due to reduced susceptibility, and does not include any change in exposure of susceptible individuals to infected herd mates. Next, prevalence decreases further because the initial decrease in prevalence reduces the exposure of susceptible individuals to infected herd mates. This additional decrease represents the indirect response to selection via the “social” environment. Without genetic variation in infectivity, the direct response makes up a fraction P of the total response in prevalence, and the indirect response a fraction $1 - P$.

The feedback mechanisms outlined here will also play a role in macroparasitic infections. For example, also for macroparasitic infections infectivity will have a nonlinear effect, susceptibility does not have to be zero to eradicate an infection, and prevalence will go down more than linear with a genetic decrease in infectivity. However, we did not investigate how this works out precisely, for example for gastro intestinal parasites.

Herd immunity

In Figure 8A, the infection ultimately goes extinct due to mass selection for individual infection status. This happens due to a phenomenon known as herd immunity (Fine 1993). In the final generation, the infection disappears because R_0 falls below a value of one; not because all the individuals have become fully resistant to infection. This result is similar to the eradication of an infection by means of vaccination, which also does not require full immunity of all individuals and can also be achieved when only part of a population is vaccinated (Anderson and May 1985). As can be seen in Figure 10 and in simulation results of Hulst et al. (2021), herd immunity develops over cycles of the transmission-recovery loop. Thus, the full benefits of genetic selection or vaccination do not manifest immediately, as it takes some time for a population to converge to the new endemic steady state.

The relevance of herd immunity for response to genetic selection can be illustrated using the data underlying Figure 8A. For the population starting at a prevalence of 0.5, the contact rate is equal to two, and the mean breeding value for log-susceptibility is equal to zero in the initial generation ($c = 2$, $\bar{A}_{i0} = 0$, so that $R_0 \approx ce^{\bar{A}_{i0}} = 2$). In the final generation, the mean breeding value for log-susceptibility has dropped to -0.73 , so that $R_0 \approx 2e^{-0.73} = 0.96$. Hence, $R_0 < 1$, explaining extinction. However, if the average individual of the final generation would have been exposed to the infection pressure of the first generation, then the expected prevalence for this individual would have been 0.32 (from Equation 20a, with $R_{0,i} = 0.96$ and $P = 0.5$). Hence, the individual would have been infected 32% of the time. Nevertheless, in a population consisting entirely of this type of individual, as is the case in the final generation, the infection will no longer be present in the long term. This example illustrates the relevance of reduced exposure due to indirect effects for herd immunity and for response to selection of infectious diseases.

Relationship to previous work

Bishop and co-workers have pioneered the integration of quantitative genetics and epidemiology for livestock populations (see

references in the Introduction). Some of their work considers both the prevalence of an infection and the negative effect of the infection on performance traits (resilience) in an integrated approach, mostly using stochastic simulation. In this study, in contrast, we focus exclusively on prevalence, since our primary purpose was to develop a quantitative genetic theory for the endemic prevalence of an infection. In particular, we aimed to find expressions for the breeding value and the additive genetic variance in the endemic prevalence. Our results show that these are fundamentally different from quantitative genetic expressions for noncommunicable traits, exhibiting a very large component due to IGE. The effect of an infectious disease on performance traits, in contrast, can be modelled using classical quantitative genetic approaches, such as reaction norm models where trait values are regressed on pathogen load. Hence, resilience may not exhibit indirect genetic variation, i.e., when it is independent of susceptibility, recovery or infectivity, and there is no need to include resilience in theoretical models for prevalence.

MacKenzie and Bishop (2001) and Tsairidou et al. (2019) investigated the prediction of response to selection in the prevalence of infectious diseases, considering both quantitative genetics and epidemiology. MacKenzie and Bishop (2001) directly modeled a constant rate of genetic improvement for the transmission parameter β , treated as a genetic property of the susceptible (i.e., recipient) individual only, and used a stochastic epidemic model to study the impact of genetic improvement in β on R_0 and on the probability of a major epidemic. Tsairidou et al. (2019) used a similar approach, but considered both infectivity and susceptibility. They directly modelled response in susceptibility and infectivity, assuming a fixed accuracy of selection for these two traits, and also used a stochastic epidemic model to study the impact of genetic improvement in susceptibility and infectivity on R_0 and on the severity of the epidemic. Hence, these two studies combine a classical quantitative genetic approach for response to selection for the parameters of an epidemiological model with stochastic simulation of epidemics. In this study, in contrast, we extend quantitative genetic theory to include R_0 and the endemic prevalence, aiming to understand the genetic variation and potential response to selection in these population-level parameters. Hence, we aim to bring epidemiology into the quantitative genetic domain, rather than to combine classical quantitative genetics models for epidemiological parameters with simulation of epidemics.

The breeding value and additive genetic variance for the logarithm of R_0 are central to this work. Anche et al. (2014) and Biemans et al. (2019) presented a breeding value and additive genetic variance for R_0 , rather than its logarithm. Anche et al. (2014) considered a two locus model with additive effects for susceptibility and infectivity. They derived a breeding value for R_0 using partial derivatives of R_0 with respect to the allele frequencies at each of the two loci. While their model is additive for susceptibility and infectivity, it contains some epistasis for R_0 because R_0 depends on the product of these two parameters. For locus-based models with a few loci of fixed effect, it is probably not very relevant on which scale the model is additive (if any). For polygenic models, in contrast, an additive model on the scale of susceptibility and infectivity, or on the scale of R_0 , may result in negative values for R_0 and in an unrealistically large additive genetic variance in R_0 with recurrent genetic selection for lower prevalence. Hence, for polygenic traits, an additive model on the log scale is more appropriate, as argued above.

Biemans et al. (2019) presented an expression for the additive genetic variance for R_0 treated as a polygenic traits, with the aim

to quantify the amount of heritable variation in R_0 in a data analysis. Their expression extends the concept of Anche et al. (2014) to the polygenic case, but can also be interpreted as the variance of a first-order Taylor-series linearization of R_0 , assuming independence of susceptibility and infectivity (combining Equations 4 and 5 of the current manuscript). The expression of Biemans et al. (2019) is suitable when the objective is to find a point estimate for the additive genetic variance in R_0 in a population, and when the additive genetic variances in susceptibility and infectivity are not too large and susceptibility and infectivity are independent. For a quantitative genetic theory of R_0 , however, an approach based on the breeding value for the logarithm of R_0 is superior, as argued above.

Utilization of hidden genetic variation for genetic improvement

In this study, we have shown that a fraction $1 - P$ of the full individual genetic effect on the endemic prevalence represents an IGE, because only a fraction P of the full effect surfaces in the infection status of the individual itself (excluding genetic variation in infectivity; Equation 33 and Appendix E). In other words, a fraction $1 - P$ of the individual genetic effects of susceptibility and recovery on the prevalence are hidden to direct selection and classical genetic analysis. Nevertheless, results in Figure 8 show that prevalence responds rapidly to selection, particularly when prevalence is small. Hence, prevalence responds faster to selection when a greater proportion of its heritable variation is hidden, and when heritability is low (Figure 7A), which seems a paradox.

However, the IGEs due to susceptibility and recovery are a special kind, because they are fully correlated to the corresponding DGE. For each of the two traits, there is only a single genetic effect (A_{I_y} and A_{I_x} , respectively), which has both a direct effect and an indirect effect on the prevalence. Hence, when selection changes the mean DGE, the mean IGE changes correspondingly. This can be seen from Equation (38b), where the term σ_{A_y}/P represents the full additive genetic standard deviation in prevalence (as is clear from Equation 33), while the accuracy ($\rho_{A_y, \bar{y}}$) refers to selection for the direct effect only. Hence, without genetic variation in infectivity, the total response of prevalence to selection, either for individual infection status or for any other selection criterion, can be interpreted as the sum of a direct response in DGE and a correlated response in IGE,

$$\begin{aligned} R_{P, \text{direct}} &= i \rho_{A_y, \bar{y}} \sigma_{A_y} \\ R_{P, \text{correlated}} &= i \rho_{A_y, \bar{y}} \sigma_{A_y} \frac{1 - P}{P} \end{aligned}$$

and the sum of $R_{P, \text{direct}}$ and $R_{P, \text{correlated}}$ is equal to Equation (38b). The $R_{P, \text{direct}}$ is the expected response to selection based on ordinary genetic analysis and estimated breeding values for individual infection status. The $R_{P, \text{correlated}}$ represents the additional response due to IGE. The direct response occurs immediately in the first cycle of the transmission loop (Figure 10B), while the indirect response manifests gradually over several cycles of the transmission-recovery loop, particularly when prevalence is small (see also result in Hulst et al. 2021).

The response due to the IGE of susceptibility and recovery arises naturally when selecting for lower individual infection status (i.e., for the direct effect); it does not require any specific measures of the breeder. Thus, on the one hand, our results imply that response to genetic selection against infectious diseases should be considerably greater than currently believed, even when no changes are made to the selection strategy. While

empirical studies are scarce, the available results support this expectation (discussed in Hulst et al. 2021).

On the other hand, however, classical selection for direct effects is not the optimal way to reduce prevalence, for the following two reasons. First, classical selection does not target genetic effects on infectivity, because an individual's infectivity does not affect its own infection status (Lipschutz-Powell et al. 2012). Hence, infectivity changes merely due to a potential genetic correlation with susceptibility and/or recovery. When this correlation is unfavorable, infectivity will increase and response in prevalence will be smaller than expected based on the genetic selection differentials for susceptibility and recovery (and thus smaller than the result of Equation 38b). In theory, this could even lead to a negative net response (Griffing 1967). This is similar to the case with social behavior-related IGEs on survival in laying hens and Japanese quail, where selection for individual survival has sometimes increased mortality (Craig and Muir 1996; Muir 2005). This scenario seems unlikely for infectious diseases, but at present we lack knowledge of the multivariate genetic parameters of susceptibility, infectivity and recovery to make well-founded statements.

Second, even in the absence of genetic variation in infectivity, individual selection for susceptibility and recovery is nonoptimal because the accuracy of selection is limited due to limited heritability, particularly at low prevalence (Figure 7A). The response to selection in traits affected by IGE can be increased by using kin selection and/or group selection (Griffing 1976; Muir 1996; Bijma 2011), and by including IGE in the genetic analysis (Muir 2005, Bijma et al. 2007b; Muir et al. 2013; Anacleto et al. 2015; Biemans et al. 2019; Pooley et al. 2020). Kin selection occurs when transmission takes place between related individuals, for example within groups of relatives (Anche et al. 2014). Group selection refers to the selection of parents for the next generation based on the prevalence in the group in which transmission takes place, rather than on individual infection status (Griffing 1976). Both theoretical and empirical work shows that kin and group selection lead to utilization of the full genetic variation, including both DGE and IGE (Griffing 1976; Muir 1996, 2005; Bijma and Wade 2008; Bijma 2010, 2011). For infectious diseases, the work of Anche et al. (2014) illustrates the effect of kin selection, where favorable alleles for susceptibility increase much faster in frequency when disease transmission is between related individuals. Simulation studies on IGE in pig populations suggest that the benefits of kin selection also apply to breeding schemes based on genomic prediction (Chu et al. 2021).

Do pathogens create kin selection?

Exposure to infectious pathogens is a major driver of the evolution of host populations by natural selection, both in animals and plants (reviewed in Karlsson et al. 2014 and Ebert and Fields 2020). In the human species, for example, a study of genetic variation in 50 worldwide populations reveals that exposure to infectious pathogens is the primary driver of local adaptation and the strongest selective force that shapes the human genome (Barreiro and Quintana-Murci 2010; Fumagalli et al. 2011). The key role of infectious pathogens in natural selection, together with the large contribution of IGE to the genetic variation in prevalence in the host population, indicates that IGE must have been an important fitness component in the evolutionary history of populations. This, in turn, suggests that associating with kin may have evolved as an adaptive behavior. In other words, the key role of infectious diseases in natural selection might lead to social structures where individuals associate preferably with kin,

because such behavior has indirect fitness benefits. This is because interactions among kin lead to utilization of the full heritable variation in fitness, including both DGE and IGE (Bijma 2010), and thus considerably accelerate response of fitness to selection. At low to moderate levels of the endemic prevalence the indirect genetic variance in prevalence might be sufficiently large for such behavior to evolve, even in the absence of direct fitness benefits such as preferential behavior toward kin. While this is a complex issue requiring careful quantitative modelling, including migration and the emergence of selfish mutants, the key role of pathogens in natural selection together with the large IGE demonstrated here strongly suggest the importance of kin selection in the history of life.

In agriculture, the implementation of kin selection may be feasible when animals can be kept in kin groups or plants can be grown in plots of a single genotype or a family in the breeding population. In many cases, however, this will not be feasible, and other methods will be required to optimally capture the IGE underlying the prevalence of infectious diseases. In particular, we need more and better phenotypic data on disease traits (Bishop and Woolliams 2014). Current developments in sensing technology and artificial intelligence enable the development of tools for large scale automated collection of longitudinal data on individual infection status, and also on the contact structure between individuals (relevant mainly in animals). These advances, together with genomic prediction and recently developed statistical methods for the estimation of the direct and IGE underlying the transmission of the infection (Biemans *et al.* 2019; Pooley *et al.* 2020) could represent a much-needed breakthrough in the artificial selection against infectious diseases in agriculture. Our results on genetic variation and response to selection suggest that such selection is way more promising than currently believed.

Data availability

An R-code to numerically find the endemic equilibrium prevalence is provided in the file [Supplementary Material 1](#) - Numerical Solution Prevalence Heterogeneity.

[Supplementary material](#) is available at GENETICS online.

Acknowledgments

The authors thank Jack C. M. Dekkers for helpful comments on the manuscript.

Funding

Funding for this study was received from the host institutions of the authors.

Conflicts of interest

The authors declare that there is no conflict of interest.

Literature cited

- Anacleto O, Cabaleiro S, Villanueva B, Saura M, Houston RD, *et al.* 2019. Genetic differences in host infectivity affect disease spread and survival in epidemics. *Sci Rep.* 9:1–12.
- Anacleto O, Garcia-Cortés LA, Lipschutz-Powell D, Woolliams JA, Doeschl-Wilson AB. 2015. A novel statistical model to estimate host genetic effects affecting disease transmission. *Genetics.* 201: 871–884.
- Anche MT, De Jong MCM, Bijma P. 2014. On the definition and utilization of heritable variation among hosts in reproduction ratio R_0 for infectious diseases. *Heredity.* 113:364–374.
- Anderson R, May R. 1979. Population biology of infectious diseases: Part I. *Nature.* 280:361–367. doi:10.1038/280361a0.
- Anderson RM, May RM. 1985. Vaccination and herd immunity to infectious diseases. *Nature.* 318:323–329.
- Aznar I, Frankena K, More SJ, O’Keeffe J, McGrath G, *et al.* 2018. Quantification of *Mycobacterium bovis* transmission in a badger vaccine field trial. *Prev Vet Med.* 149:29–37.
- Barreiro LB, Quintana-Murci L. 2010. From evolutionary genetics to human immunology: how selection shapes host defence genes. *Nat Rev Genet.* 11:17–30.
- Biemans F, de Jong MC, Bijma P. 2017. A model to estimate effects of SNPs on host susceptibility and infectivity for an endemic infectious disease. *Genet Sel Evol.* 49:1–13.
- Biemans F, De Jong MC, Bijma P. 2019. Genetic parameters and genomic breeding values for digital dermatitis in Holstein Friesian dairy cattle: host susceptibility, infectivity and the basic reproduction ratio. *Genet Sel Evol.* 51:1–13.
- Bijma P. 2010. Fisher’s fundamental theorem of inclusive fitness and the change in fitness due to natural selection when conspecifics interact. *J Evol Biol.* 23:194–206.
- Bijma P. 2011. A general definition of the heritable variation that determines the potential of a population to respond to selection. *Genetics.* 189:1347–1359.
- Bijma P. 2020. The Price equation as a bridge between animal breeding and evolutionary biology. *Philos Trans R Soc Lond B Biol Sci.* 375:20190360.
- Bijma P, Muir WM, Van Arendonk JA. 2007. Multilevel selection 1: quantitative genetics of inheritance and response to selection. *Genetics.* 175:277–288.
- Bijma P, Wade MJ. 2008. The joint effects of kin, multilevel selection and indirect genetic effects on response to genetic selection. *J Evol Biol.* 21:1175–1188.
- Bishop S, Doeschl-Wilson AB, Woolliams JA. 2012. Uses and implications of field disease data for livestock genomic and genetics studies. *Front Genet.* 3:114.
- Bishop SC, MacKenzie KM. 2003. Genetic management strategies for controlling infectious diseases in livestock populations. *Genet Sel Evol.* 35:S3–S17.
- Bishop SC, Stear MJ. 1997. Modelling responses to selection for resistance to gastro-intestinal parasites in sheep. *Anim Sci.* 64: 469–478.
- Bishop SC, Stear MJ. 1999. Genetic and epidemiological relationships between productivity and disease resistance: gastro-intestinal parasite infection in growing lambs. *Anim Sci.* 69:515–524.
- Bishop SC, Stear MJ. 2003. Modelling of host genetics and resistance to infectious diseases: understanding and controlling nematode infections. *Vet Parasitol.* 115:147–166.
- Bishop SC, Woolliams JA. 2010. On the genetic interpretation of disease data. *PLoS One.* 5:e8940. doi:10.1371/journal.pone.0008940.
- Bishop SC, Woolliams JA. 2014. Genomics and disease resistance studies in livestock. *Livest Sci.* 166:190–198.
- Bolker BM, Grenfell BT. 1996. Impact of vaccination on the spatial correlation and persistence of measles dynamics. *Proc Natl Acad Sci USA.* 93:12648–12653.
- Chu TT, Henryon M, Jensen J, Ask B, Christensen OF. 2021. Statistical model and testing designs to increase response to selection with constrained inbreeding in genomic breeding programs for pigs affected by social genetic effects. *Genet Sel Evol.* 53:1–16.

- Craig JV, Muir WM. 1996. Group selection for adaptation to multiple-hen cages: beak-related mortality, feathering, and body weight responses. *Poult Sci.* 75:294–302.
- De Villemereuil P, Schielzeth H, Nakagawa S, Morrissey M. 2016. General methods for evolutionary quantitative genetic inference from generalized mixed models. *Genetics.* 204:1281–1294.
- DeJong MCM, vanderPoel WHM, Kramps JA, Brand A, vanOirschot JT. 1996. Quantitative investigation of population persistence and recurrent outbreaks of bovine respiratory syncytial virus on dairy farms. *Am J Vet Res.* 57:628–633.
- Dempster ER, Lerner IM. 1950. Heritability of threshold characters. *Genetics.* 35:212–236.
- Diekmann O, Heesterbeek H, Britton T. 2012. *Mathematical Tools for Understanding Infectious Disease Dynamics*, Vol. 7. Princeton: Princeton University Press.
- Diekmann O, Heesterbeek JAP, Metz JAJ. 1990. On the definition and the computation of the basic reproduction ratio in models for infectious diseases in heterogeneous populations. *J Math Biol.* 28:365–382. doi:10.1007/BF00178324.
- Doeschl-Wilson AB, Davidson R, Conington J, Roughsedge T, Hutchings MR, Villanueva B. 2011. Implications of host genetic variation on the risk and prevalence of infectious diseases transmitted through the environment. *Genetics.* 188:683–693.
- Doeschl-Wilson A, Knap PW, Opriessnig T, More SJ. 2021. Livestock disease resilience: from individual to herd level. *Animal*. doi:10.1016/j.animal.2021.100286.
- Ebert D, Fields PD. 2020. Host–parasite co-evolution and its genomic signature. *Nat Rev Genet.* 21:754–768.
- EFSA Panel on Animal Health and Welfare (AHAW). 2012. Scientific Opinion on Review of the European Union Summary Report on trends and sources of zoonoses, zoonotic agents and food-borne outbreaks in 2009 and 2010 specifically for the data related to bovine tuberculosis, *Echinococcus*, Q fever, brucellosis and non-food borne diseases. *EFSA J.* 10:2765.
- Fine PE. 1993. Herd immunity: history, theory, practice. *Epidemiol Rev.* 15:265–302.
- Fisher RA. 1919. XV.—The correlation between relatives on the supposition of Mendelian inheritance. *Trans R Soc Edinb.* 52:399–433.
- Fumagalli M, Sironi M, Pozzoli U, Ferrer-Admetlla A, Ferrer-Admetlla A, et al. 2011. Signatures of environmental genetic adaptation pinpoint pathogens as the main selective pressure through human evolution. *PLoS Genet.* 7:e1002355.
- Gianola D. 1982. Theory and analysis of threshold characters. *J Anim Sci.* 54:1079–1096.
- Gillespie DT. 1977. Exact stochastic simulation of coupled chemical reactions. *J Phys Chem.* 81:2340–2361.
- Greenhalgh D, Diekmann O, de Jong MCM. 2000. Subcritical endemic steady states in mathematical models for animal infections with incomplete immunity. *Math Biosci.* 165:1–25. doi:10.1016/S0025-5564(00)00012-2.
- Griffing B. 1967. Selection in reference to biological groups I. Individual and group selection applied to populations of unordered groups. *Aust J Biol Sci.* 20:127–140.
- Griffing B. 1976. Selection in reference to biological groups. V. Analysis of full-sib groups. *Genetics.* 82:703–722.
- Griffing B. 1977. Selection for populations of interacting genotypes. In: E Pollack, O Kempthorne, TB Bailey, editors. *Proceedings of the International Congress on Quantitative Genetics*, August 16–21, 1976. Ames, IW: Iowa State University Press. p. 413–434.
- Heringstad B, Klemetsdal G, Steine T. 2007. Selection responses for disease resistance in two selection experiments with Norwegian red cows. *J Dairy Sci.* 90:2419–2426.
- Hethcote HW. 1989. Three basic epidemiological models. In: Levin SA, Hallam TG, Gross LJ, editors. *Applied Mathematical Ecology*. Berlin, Heidelberg: Springer. p. 119–144.
- Hill WG, Mulder HA. 2010. Genetic analysis of environmental variation. *Genet Res (Camb).* 92:381–395.
- Hulst AD, de Jong MC, Bijma P. 2021. Why genetic selection to reduce the prevalence of infectious diseases is way more promising than currently believed. *Genetics.* 217:iyab024.
- Karlsson EK, Kwiatkowski DP, Sabeti PC. 2014. Natural selection and infectious disease in human populations. *Nat Rev Genet.* 15:379–393.
- Keeling MJ, Rohani P. 2011. *Modeling Infectious Diseases in Humans and Animals*. Princeton University Press.
- Kermack WO, McKendrick AG. 1927. A contribution to the mathematical theory of epidemics. *Proc R Soc Lond A.* 115:700–721.
- Knap PW, Doeschl-Wilson A. 2020. Why breed disease-resilient livestock, and how? *Genet Sel Evol.* 52:1–18.
- Lipschutz-Powell D, Woolliams JA, Bijma P, Doeschl-Wilson AB. 2012. Indirect genetic effects and the spread of infectious disease: are we capturing the full heritable variation underlying disease prevalence? *PLoS One.* 7:e39551.
- Lipschutz-Powell D, Woolliams JA, Doeschl-Wilson AB. 2014. A unifying theory for genetic epidemiological analysis of binary disease data. *Genet Sel Evol.* 46:15–12.
- MacKenzie K, Bishop SC. 1999. A discrete-time epidemiological model to quantify selection for disease resistance. *Anim Sci.* 69:543–551.
- MacKenzie K, Bishop SC. 2001. Utilizing stochastic genetic epidemiological models to quantify the impact of selection for resistance to infectious diseases in domestic livestock. *J Anim Sci.* 79:2057–2065.
- Martin P, Barkema H, Brito L, Narayana S, Miglior F. 2018. Symposium review: novel strategies to genetically improve mastitis resistance in dairy cattle. *J Dairy Sci.* 101:2724–2736.
- McCullagh P, Nelder JA. 2019. *Generalized Linear Models*. Boca Raton, Florida: Routledge.
- Metz JAJ. 1978. The epidemic in a closed population with all susceptibles equally vulnerable; some results for large susceptible populations and small initial infections. *Acta Biotheor.* 27:75–123.
- Moore AJ, Brodie ED, III, Wolf JB. 1997. Interacting phenotypes and the evolutionary process: I. Direct and indirect genetic effects of social interactions. *Evolution.* 51:1352–1362.
- Muir WM. 1996. Group selection for adaptation to multiple-hen cages: selection program and direct responses. *Poult Sci.* 75:447–458.
- Muir WM. 2005. Incorporation of competitive effects in forest tree or animal breeding programs. *Genetics.* 170:1247–1259.
- Muir WM, Bijma P, Schinckel A. 2013. Multilevel selection with kin and non-kin groups, experimental results with Japanese quail (*Coturnix japonica*). *Evolution.* 67:1598–1606.
- Nelder JA, Wedderburn RW. 1972. Generalized linear models. *J R Stat Soc A.* 135:370–384.
- Nieuwhof GJ, Conington J, Bishop SC. 2009. A genetic epidemiological model to describe resistance to an endemic bacterial disease in livestock: application to footrot in sheep. *Genet Sel Evol.* 41:19.
- Pooley CM, Marion G, Bishop SC, Bailey RI, Doeschl-Wilson AB. 2020. Estimating individuals' genetic and non-genetic effects underlying infectious disease transmission from temporal epidemic data. *PLoS Comput Biol.* 16:e1008447.
- Robertson A. 1950. Proof that the additive heritability on the p scale is given by the expression z_{2h}^2/pq . *Genetics.* 35:234–236.
- Rushton J. 2009. *The Economics of Animal Health and Production*. Wallingford: CABI.

- Russell GE. 2013. Plant Breeding for Pest and Disease Resistance: Studies in the Agricultural and Food Sciences. London, UK: Butterworths.
- SanCristobal-Gaudy M, Elsen JM, Bodin L, Chevalet C. 1998. Prediction of the response to a selection for canalisation of a continuous trait in animal breeding. *Genet Sel Evol.* 30:423–451.
- Schrag SJ, Wiener P. 1995. Emerging infectious disease: what are the relative roles of ecology and evolution? *Trends Ecol Evol.* 10: 319–324.
- Springbett AJ, MacKenzie K, Woolliams JA, Bishop SC. 2003. The contribution of genetic diversity to the spread of infectious diseases in livestock populations. *Genetics.* 165:1465–1474. doi:10.1093/genetics/165.3.1465.
- Thanner S, Drissner D, Walsh F. 2016. Antimicrobial resistance in agriculture. *mBio.* 7:e02227–15.
- Tsairidou S, Anacleto O, Woolliams JA, Doeschl-Wilson A. 2019. Enhancing genetic disease control by selecting for lower host infectivity and susceptibility. *Heredity.* 122:742–758.
- Walsh B, Lynch M. 2018. Evolution and Selection of Quantitative Traits. Oxford, UK: Oxford University Press.
- Weiss GH, Dishon M. 1971. On the asymptotic behavior of the stochastic and deterministic models of an epidemic. *Math Biosci.* 11: 261–265.
- Wolf JB, Brodie ED, III, Cheverud JM, Moore AJ, Wade MJ. 1998. Evolutionary consequences of indirect genetic effects. *Trends Ecol Evol.* 13:64–69.
- Wright AJ. 1985. Selection for improved yield in inter-specific mixtures or intercrops. *Theor Appl Genet.* 69:399–407.

Communicating editor: M. Goddard

Appendices

Appendix A

R_0 with heterogeneity and log-normally distributed susceptibility, infectivity en recovery

We assume that the transmission rate from infected individual j to susceptible individual i is proportional to the product of the infectivity of j and the susceptibility of i (Equation 5),

$$\beta_{ij} = c\gamma_i\varphi_j.$$

So there is no interaction between i and j . (This property is known as separable mixing in the epidemiological literature; Diekmann et al. 1990, 2013). Moreover, we assume that susceptibility, infectivity, and recovery follow a log-normal distribution (Equations 6 and 7). We also assume that the population is not very small, so that in the early phase of an endemic where only few individuals are infected, the composition of the remaining susceptible individuals is not affected.

Because R_0 refers to the “total number of individuals that become infected by a typical infected individual over its entire infectious lifetime,” we define an individual lifetime infectivity, which is the product of an individual’s infectivity per unit of time and the average duration of its infectious lifetime,

$$\phi_i = \varphi_i/\alpha_i,$$

which follows a log-normal distribution with parameters following from those of φ and α . Hence, we have condensed our three genetic effects into two.

We can find R_0 from

$$R_0 = c \bar{\gamma} \phi_{\text{typ}}$$

where ϕ_{typ} is the lifetime infectivity of the typical infected individual, and $\bar{\gamma}$ is the simple average susceptibility in the population,

$$\bar{\gamma} = \int_0^\infty \gamma g(\gamma) d\gamma.$$

where $g(\gamma)$ is the pdf of γ . We can use the simple average of susceptibility in this expression because we assume the population is large.

With separable mixing, the typical infected individual is created immediately in the first generation of disease transmission. This is the case because there is no interaction between γ and φ , so that the properties of the typical infected individual are determined entirely by susceptibility. Hence, the pdf of γ for the typical infected individual follows from weighing $g(\gamma)$ by γ ,

$$g_{\text{typ}}(\gamma) = \frac{1}{\bar{\gamma}} \gamma g(\gamma)$$

Since the properties of the typical infected individual depend on susceptibility only, we can find ϕ_{typ} by averaging ϕ over its distribution conditional on γ , and subsequently averaging over the distribution of γ ,

$$\phi_{\text{typ}} = \int_0^\infty \left(\int_0^\infty \phi f(\phi|\gamma) d\phi \right) g_{\text{typ}}(\gamma) d\gamma$$

Hence, we now have the elements of R_0 , but still need to solve the integral expression.

Because conditional Normal distributions are also Normal and the logarithm is a bijective function, $\phi|\gamma$ follows a log-normal distribution with parameters being the conditional mean and variance of the Normal distribution,

$$\phi|\gamma \sim \text{Lnorm}\left(\mu = b_{\phi,\gamma}A_{l\gamma}; \sigma^2 = (1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2\right)$$

with $b_{\phi,\gamma} = \text{cov}(A_{l\gamma}, A_{l\phi})/\text{var}(A_{l\gamma})$ denoting the regression coefficient of $A_{l\phi}$ on $A_{l\gamma}$, and $\rho_{\gamma,\phi}^2 = \text{cov}^2(A_{l\gamma}, A_{l\phi})/[\text{var}(A_{l\gamma})\text{var}(A_{l\phi})]$ the squared correlation, where $A_{l\phi}$ denotes the breeding value for logarithm of lifetime infectivity.

Hence, the inner integral is the mean of a log-normal variate, which is of the form $\exp\left(\mu + \frac{\sigma^2}{2}\right)$,

$$\begin{aligned} \int_0^\infty \phi f(\phi|\gamma) d\phi &= E[\phi|\gamma] = \exp\left(b_{\phi,\gamma}A_{l\gamma} + \frac{1}{2}(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2\right) \\ &= e^{\frac{1}{2}(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2} e^{b_{\phi,\gamma}A_{l\gamma}}. \end{aligned}$$

Because the first term of this expression is a constant,

$$\phi_{\text{typ}} = e^{\frac{1}{2}(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2} \int_0^\infty e^{b_{\phi,\gamma}A_{l\gamma}} g_{\text{typ}}(\gamma) d\gamma.$$

Substituting $g_{\text{typ}}(\gamma) = \frac{1}{\bar{\gamma}} \gamma g(\gamma)$, and replacing $g(\gamma)$ by the corresponding log-normal pdf yields

$$\int_0^\infty e^{b_{\phi,\gamma}A_{l\gamma}} g_{\text{typ}}(\gamma) d\gamma = \frac{1}{\bar{\gamma}\sigma\sqrt{2\pi}} \int_0^\infty e^{\frac{A^2}{2\sigma^2} + bA} d\gamma$$

where we simplified the notation for brevity, using $\sigma^2 = \sigma_{A_{l\gamma}}^2$, $b = b_{\phi,\gamma}$, and $A = A_{l\gamma}$.

Next, we change variable, using $d\gamma = e^A dA$, and adjust the bounds accordingly,

$$\frac{1}{\bar{\gamma}\sigma\sqrt{2\pi}} \int_{-\infty}^\infty e^{\frac{A^2}{2\sigma^2} + bA} e^A dA = \frac{1}{\bar{\gamma}\sigma\sqrt{2\pi}} \int_{-\infty}^\infty e^{\frac{A^2}{2\sigma^2} + (1+b)A} dA$$

Solving the integral term in Mathematica-online yields

$$\begin{aligned} \frac{1}{\bar{\gamma}\sigma\sqrt{2\pi}} \int_{-\infty}^\infty e^{\frac{A^2}{2\sigma^2} + (1+b)A} dA &= \frac{1}{\bar{\gamma}} e^{\frac{1}{2}\sigma_{A_{l\gamma}}^2 (1+b_{\phi,\gamma})^2} \\ \phi_{\text{typ}} &= e^{\frac{1}{2}(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2} \frac{1}{\bar{\gamma}} e^{\frac{1}{2}\sigma_{A_{l\gamma}}^2 (1+b_{\phi,\gamma})^2} \\ \phi_{\text{typ}} &= \frac{1}{\bar{\gamma}} e^{\frac{1}{2}[(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2 + \sigma_{A_{l\gamma}}^2 (1+b_{\phi,\gamma})^2]} \\ R_0 &= c \bar{\gamma} \phi_{\text{typ}} = c e^{\frac{1}{2}[(1 - \rho_{\gamma,\phi}^2)\sigma_{A_{l\phi}}^2 + \sigma_{A_{l\gamma}}^2 (1+b_{\phi,\gamma})^2]} \end{aligned}$$

Using $\bar{\gamma} = e^{\frac{1}{2}\sigma_{A_{l\gamma}}^2}$ and $\bar{\phi} = e^{\frac{1}{2}\sigma_{A_{l\phi}}^2}$ this simplifies to Equations (14) and (16) of the main text,

$$\begin{aligned}\phi_{\text{typ}} &= \bar{\phi} e^{\sigma_{A_{ly}A_{ly}}} \\ R_0 &= c \bar{\gamma} \bar{\phi} e^{\sigma_{A_{ly}A_{ly}}}.\end{aligned}$$

Further simplification follows from expressing $\bar{\gamma}$, $\bar{\phi}$ and $e^{\sigma_{A_{ly}A_{ly}}}$ in terms of variances and covariances of γ , ϕ and α .

$$\begin{aligned}\bar{\gamma} &= e^{\frac{1}{2}\sigma_{A_{ly}}^2} \\ \phi_i &= \frac{\phi_i}{\alpha_i} = e^{(A_{ly,i} - A_{ly,i})} = e^{A_{ly,i}}\end{aligned}$$

where $A_{ly,i} = A_{ly,i} - A_{ly,i}$, which is the breeding value for the logarithm of lifetime infectivity, with

$$\text{var}(A_{ly,i}) = \sigma_{A_{ly}}^2 - 2\sigma_{A_{ly}A_{ly}} + \sigma_{A_{ly}}^2$$

From the properties of the log-normal distribution,

$$\bar{\phi} = e^{\frac{1}{2}\sigma_{A_{ly}}^2} = e^{\frac{1}{2}(\sigma_{A_{ly}}^2 - 2\sigma_{A_{ly}A_{ly}} + \sigma_{A_{ly}}^2)}$$

Furthermore,

$$e^{\sigma_{A_{ly}A_{ly}}} = e^{(\sigma_{A_{ly}A_{ly}} - \sigma_{A_{ly}A_{ly}})}$$

Substitution of the expressions for $\bar{\gamma}$, $\bar{\phi}$, and $e^{\sigma_{A_{ly}A_{ly}}}$ into $R_0 = c \bar{\gamma} \bar{\phi} e^{\sigma_{A_{ly}A_{ly}}}$ yields

$$R_0 = c e^{\frac{1}{2}\sigma_{A_{ly}}^2} e^{\frac{1}{2}(\sigma_{A_{ly}}^2 - 2\sigma_{A_{ly}A_{ly}} + \sigma_{A_{ly}}^2)} e^{(\sigma_{A_{ly}A_{ly}} - \sigma_{A_{ly}A_{ly}})}$$

$$R_0 = c e^{\frac{1}{2}(\sigma_{A_{ly}}^2 + \sigma_{A_{ly}}^2 + \sigma_{A_{ly}}^2 + 2\sigma_{A_{ly}A_{ly}} - 2\sigma_{A_{ly}A_{ly}} - 2\sigma_{A_{ly}A_{ly}})} R_0 = c e^{\frac{1}{2}\sigma_{A_{ly}}^2}.$$

The right-hand side of this expression is identical to the mean genotypic value for R_0 (Equation 12).

Appendix B

Numerical solution to find the endemic prevalence with heterogeneity

To find the endemic prevalence, P , we partition the population into types, i , and numerically solve the expressions

$$P_i = \frac{\mathcal{R}_{0,i}P}{\mathcal{R}_{0,i}P + 1}$$

and

$$\mathcal{R}_{0,i} = \frac{c\gamma_i\bar{\phi}_{\text{inf}}}{\alpha_i}$$

for P . Here, we develop this numerical solution for the case where susceptibility, infectivity and recovery follow a log-normal distribution, assuming separable mixing (see Appendix A).

As can be seen from the expression for $\mathcal{R}_{0,i}$, the endemic prevalence for a type depends on both its susceptibility (γ_i) and its recovery rate (α_i). Individuals with above-average susceptibility are over-represented among the infecteds in the endemic equilibrium, whereas individuals with above-average recovery are under-represented. Hence, as can be seen from the expression for

$\mathcal{R}_{0,i}$, the partitioning into types should be based on γ_i/α_i . Therefore, we define

$$\begin{aligned}\theta_i &= \frac{\gamma_i}{\alpha_i} = e^{A_{\theta,i}} \\ A_{\theta,i} &= A_{\gamma,i} - A_{\alpha,i} \\ \theta &\sim \text{Lnorm}(\mu_{A_{\theta}} = 0, \sigma_{A_{\theta}}^2 = \sigma_{A_{\gamma}}^2 - 2\sigma_{A_{\gamma}A_{\alpha}} + \sigma_{A_{\alpha}}^2)\end{aligned}$$

To numerically solve the two equations given above, we also need $\bar{\phi}_{\text{inf}}$. The $\bar{\phi}_{\text{inf}}$ will depend on the θ_i of the infecteds when infectivity is correlated to susceptibility and/or recovery. Hence, we need the distribution of $\phi|\theta$, which follows from

$$\begin{aligned}\sigma_{A_{\theta}A_{\phi}} &= \sigma_{A_{\gamma}A_{\phi}} - \sigma_{A_{\alpha}A_{\phi}} \\ b_{\phi\theta} &= \frac{\sigma_{A_{\theta}A_{\phi}}}{\sigma_{A_{\theta}}^2} \\ \rho_{\phi\theta} &= \sigma_{A_{\theta}A_{\phi}} / \sigma_{A_{\theta}} \sigma_{A_{\phi}} \\ E(A_{\phi}|A_{\theta}) &= b_{\phi\theta}A_{\theta} = \mu_{\phi|\theta} \\ \text{var}(A_{\phi}|A_{\theta}) &= (1 - \rho_{\phi\theta}^2)\sigma_{A_{\phi}}^2 = \sigma_{\phi|\theta}^2\end{aligned}$$

so that

$$\phi|\theta \sim \text{Lnorm}(\mu_{\phi|\theta} = b_{\phi\theta}A_{\theta}, \sigma_{\phi|\theta}^2 = (1 - \rho_{\phi\theta}^2)\sigma_{A_{\phi}}^2)$$

From the log-normal distribution:

$$E(\phi|\theta) = e^{\mu_{\phi|\theta} + \frac{1}{2}\sigma_{\phi|\theta}^2}$$

Hence, we can partition θ into classes i , with

$$\begin{aligned}\mathcal{R}_{0,i} &= c \theta_i \bar{\phi}_{\text{inf}} \\ \bar{\phi}_{\text{inf}} &= \frac{1}{P} \sum_i f_i P_i E(\phi_i|\theta_i) \\ P_i &= \frac{\mathcal{R}_{0,i}P}{\mathcal{R}_{0,i}P + 1} \\ P &= \sum_i f_i P_i\end{aligned}$$

where f_i is the fraction of individuals of type i , $f_i = N_i/N$, and P_i is the prevalence in type i , $P_i = I_i/N_i$. The numerical solution follows from iterating on these four equations. An R-code is in [Supplementary Material 1](#).

Appendix C

Methods for simulation of epidemics and validation of prevalence and genotypic value for individual disease status

We simulated epidemics according to standard epidemiological theory to validate the numerical solution of the endemic prevalence (Equations 20a and 20b) and the genotypic values for binary disease status (Equation 24). We considered two compartments of individuals, susceptible individuals (S) and infected (I) individuals, and a so-called stochastic SIS-model where susceptible individuals can become infected, and infected individuals can recover and then immediately become susceptible again (Weiss and Dishon 1971). For simplicity, we simulated genetic variation in susceptibility only, with $\gamma_i \sim \text{Lnorm}(0, \sigma_{A_{\gamma}}^2)$.

To limit Monte-Carlo error, we simulated a relatively large population of $N = 2,000$ genetically unrelated individuals for a total of 300,000 events (infection or recovery). We used a burn-in of

100,000 events before recording data on individual binary disease status. Hence, in the recorded data, the average individual experienced 100 events (50 infections and 50 recoveries).

The endemic was started by infecting a proportion $P_0 = 1-1/c$ of the individuals, chosen at random. Subsequently, we sampled events (infection or recovery) and the individual involved using Gillespie's algorithm (Gillespie 1977). For each infected individual, the probability of recovery was proportional to the recovery rate, α . For susceptible individual i the probability of infection was proportional to $c\gamma_i I/N$, I/N denoting the fraction of the population that is infected. Probabilities were accumulated over all individuals and scaled to a sum of 1 by dividing them by their sum. Finally, the specific event was sampled by drawing a random number, say x , from a standard uniform distribution and finding the event and the corresponding individual belonging to the probability interval $[x_l, x_h]$, where $x_l < x < x_h$. The disease status of that individual and I were updated before sampling the next event. The time of each event was not simulated. After 300,000 events, prevalence was calculated as the disease status averaged over the entire population, and also by individual, discarding the burn-in period. The regression coefficient of average individual disease status on G_y was also estimated.

Additive genetic variance in log-susceptibility was $\sigma_{A_{\gamma}}^2 = 0.3^3$. Three scenarios were considered, differing in contact rate: $c = 1.22$ giving $P = 0.2$, $c = 2$ giving $P = 5$ and $c = 5.15$ giving $P = 0.8$. Those combinations of $\sigma_{A_{\gamma}}^2$, c and P were found by numerically solving Equations (20a) and (20b). The actual prevalences observed in the simulations were equal to these numerical solutions.

Appendix D

Additive genetic variance in log-normal traits.

We assumed log-normally distributed genotypic values for susceptibility, infectivity and recovery, also resulting in a log-normal distribution for G_{R_0} and for $1 - G_P$. Hence, genetic effects are additive on the log-scale, but taking the exponent introduces some nonadditive genetic variance on the actual scale. Here, we derive the fraction of the variance that is additive on the actual scale.

Because all genetic effects had a mean of zero on the log-scale, the problem is equivalent to finding the fraction of additive variance in $z = e^x$, where $x \sim N(\mu = 0, \sigma^2)$. From the properties of the log-normal distribution, $E(z) = e^{\mu} e^{\sigma^2/2}$. With a small change $d\mu$, the mean of z becomes $e^{d\mu} e^{\sigma^2/2}$. Hence, the mean of z changes by an amount $e^{\frac{\sigma^2}{2}}(e^{d\mu} - 1)$. Since $\lim_{d\mu \rightarrow 0} e^{d\mu} = 1 + d\mu$, this change corresponds to $e^{\frac{\sigma^2}{2}} d\mu$. Thus the least-squares linear regression coefficient of z on x equals

$$b_{z,x} = e^{\sigma^2/2}.$$

For example, the linear regression coefficient of the genotypic value for prevalence (Equation 26b) on the breeding value for the logarithm of R_0 equals $b_{G_P, A_{R_0}} = c^{-1} \exp\left(\frac{1}{2} \sigma_{A_{R_0}}^2\right)$. Thus, the additive effect for z equals

$$A_z = e^{\sigma^2/2} x,$$

and additive variance in z equals

$$\sigma_{A_z}^2 = \sigma^2 e^{\sigma^2}.$$

The total variance in z follows from the properties of the log-normal distribution,

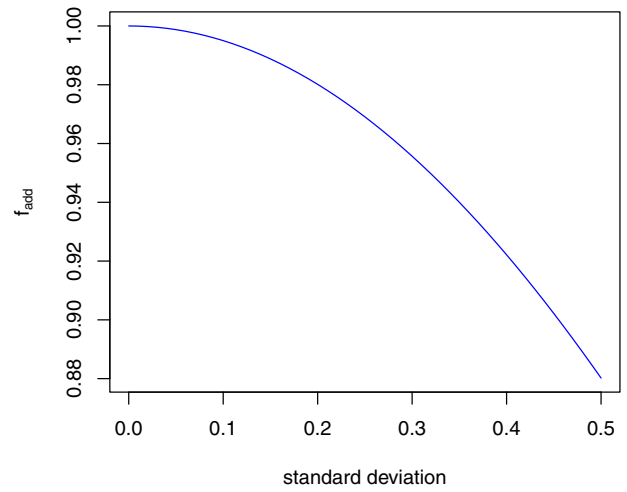


Figure D1: The additive fraction of the variance in traits following a log-normal distribution. sdx denotes the standard deviation on the normal scale.

$$\sigma_z^2 = (e^{\sigma^2} - 1)e^{\sigma^2}.$$

The additive fraction of the variance in z , therefore, equals

$$f_{\sigma_z^2} = \sigma^2 / (e^{\sigma^2} - 1)$$

Figure D1 illustrates that the additive fraction of the variance in z approaches 1 when σ^2 goes to zero. For $\sigma^2 = 0.5^2$, ~88% of the variance in z is additive. Variances on the log scale larger than 0.5^2 are unrealistic (see main text). This indicates that at least 88% of the genetic variance in susceptibility, infectivity, recovery, R_0 and prevalence is additive when they follow a log-normal distribution.

Appendix E

Breeding value for individual disease status vs breeding value for prevalence, without genetic variation in infectivity

Without genetic variation in infectivity, we have $\varphi_i = \varphi = 1$, because the scale is included in the effective contact rate c . From Equation (24), the genotypic value for individual binary disease status is,

$$G_{y,i} = \frac{R_{0,i}P}{R_{0,i}P + 1}$$

where, from Equation (20b),

$$R_{0,i} = c\gamma_i/\alpha_i.$$

From Equation (26a), the genotypic value for prevalence is,

$$G_{P,i} = 1 - \frac{1}{G_{R_0,i}}$$

where, from Equation (8),

$$G_{R_0,i} = c\gamma_i/\alpha_i.$$

Hence, without genetic variation in infectivity, $\mathcal{R}_{0,i}$ and $G_{\mathcal{R}_{0,i}}$ are identical, and we will use the symbol $G_{\mathcal{R}_{0,i}}$ in the following.

The linear approximation of the relationship between G_y and G_p follows from a comparison of their first derivatives with respect to $G_{\mathcal{R}_0}$,

$$\begin{aligned}\frac{dG_p}{dG_{\mathcal{R}_0}} &= \frac{1}{G_{\mathcal{R}_0}^2} \\ \frac{dG_y}{dG_{\mathcal{R}_0}} &= \frac{P(G_{\mathcal{R}_0}P + 1) - G_{\mathcal{R}_0}P^2}{(G_{\mathcal{R}_0}P + 1)^2} \\ &= \frac{P}{(G_{\mathcal{R}_0}P + 1)^2}\end{aligned}$$

Substituting Equation (3), assuming limited heterogeneity, yields

$$\frac{dG_y}{dG_{\mathcal{R}_0}} = \frac{G_{\mathcal{R}_0} - 1}{G_{\mathcal{R}_0}^3}.$$

Hence,

$$\frac{dG_y}{dG_{\mathcal{R}_0}} / \frac{dG_p}{dG_{\mathcal{R}_0}} = \frac{G_{\mathcal{R}_0}}{G_{\mathcal{R}_0} - 1} = P.$$

Therefore, for a small change in an individual's genotypic value for \mathcal{R}_0 , the change in its genotypic value for binary disease status is only a fraction P of the change in its genotypic value for prevalence,

$$dG_y / dG_p = P$$

Hence, when expressed relative to their mean, G_y and G_p differ approximately by a factor P (see also Figure 4 in Bijma 2020). This result is approximate, because the true relationship is nonlinear and the expression $P_{eq} = 1 - 1/\mathcal{R}_0$ is approximate with variation among individuals. For realistic magnitudes of the genetic variance, however, the nonlinearity is limited. Note that the above derivation does not require the assumption of a log-normal distribution of susceptibility and recovery.

So far, this appendix has considered genotypic values. However, the factor P also applies on the level of the breeding values, which can be shown as follows. The above derivation is based on derivatives with respect to genotypic value for \mathcal{R}_0 , say $dx/dG_{\mathcal{R}_0}$, where x is G_y or G_p . Because the relationship between P and y arises entirely via $G_{\mathcal{R}_0}$, we can translate the results to the breeding values using

$$\frac{dx}{dA_{\mathcal{R}_0}} = \frac{dx}{dG_{\mathcal{R}_0}} \frac{dG_{\mathcal{R}_0}}{dA_{\mathcal{R}_0}}.$$

The latter, $\frac{dG_{\mathcal{R}_0}}{dA_{\mathcal{R}_0}}$, is the same for x is G_y or x is G_p , so that the ratio dA_y/dA_p is the same as dG_y / dG_p . Hence, we also find

$$dA_y / dA_p = P.$$

Appendix F

Methods for observed response to selection

First a base population was generated of $N = 4000$ unrelated individuals, with genetic variation in susceptibility only. No distinction was made between males and females. For each individual, breeding values for the logarithm of susceptibility were sampled from $A_{ly} \sim N(0, 0.3^2)$, and individual susceptibility was calculated as $\gamma_i = e^{A_{ly,i}}$. The expected prevalence for the base generation was calculated as $P_0 = 1 - 1/c$, with a c of either 2 or 10, and the initial disease status of base generation individuals was sampled at random from $Bernoulli(P_0)$.

Next, an endemic was simulated following methods described in Appendix C, for a total of 15,000 events (sum of infections and recoveries), consisting of a burn-in of 10,000 events and 5,000 recorded events. The 4,000 individuals were ordered based on their mean individual disease status over the 5,000 recorded events (so based on 1.25 events on average per individual), and the 2000 individuals with the lowest values were selected as parents of the next generation (corresponding to a selected proportion of 0.5).

Selected parents were mated at random. Each pair of parents produced two offspring, resulting in $N = 4,000$ offspring. Offspring inherited the breeding value for the logarithm of susceptibility in a Mendelian fashion; $A_{ly, \text{offspring}} = \frac{1}{2}A_{ly, \text{parent1}} + \frac{1}{2}A_{ly, \text{parent2}} + N\left(0, \frac{1}{2}\sigma_{A_{ly}}^2\right)$. The initial disease status of offspring (i.e., at the start of the burn-in period of their generation) was sampled at random from $Bernoulli(P_{\text{offspring}})$, where $P_{\text{offspring}}$ denotes the expected prevalence in the offspring generation, calculated as $P_{\text{offspring}} = \max\left[1 - \frac{1}{c e^{A_{ly}}}; 0.02\right]$. The 0.02 guaranteed an average of at least 80 infected individuals at the start of the endemic in any generation, also when the expected prevalence was zero (i.e., when $1 - \frac{1}{c e^{A_{ly}}} \leq 0$). Then an endemic was started, as described above for the base generation, etc. This process was repeated until the number of infected individuals dropped to zero, implying extinction of the infection.