

Mathematical Properties of Linkage Disequilibrium Statistics Defined by Normalization of the Coefficient $D = p_{AB} - p_A p_B$

Jonathan T.L. Kang Noah A. Rosenberg

Department of Biology, Stanford University, Stanford, CA, USA

Keywords

Allele frequencies · Linkage disequilibrium · Population genetics · Statistical genetics · Statistics

Abstract

Background: Many statistics for measuring linkage disequilibrium (LD) take the form of a normalization of the LD coefficient D . Different normalizations produce statistics with different ranges, interpretations, and arguments favoring their use. **Methods:** Here, to compare the mathematical properties of these normalizations, we consider 5 of these normalized statistics, describing their upper bounds, the mean values of their maxima over the set of possible allele frequency pairs, and the size of the allele frequency regions accessible given specified values of the statistics. **Results:** We produce detailed characterizations of these properties for the statistics d and ρ , analogous to computations previously performed for r^2 . We examine the relationships among the statistics, uncovering conditions under which some of them have close connections. **Conclusion:** The results contribute insight into LD measurement, particularly the understanding of differences in the features of different LD measures when computed on the same data. © 2020 S. Karger AG, Basel

Introduction

Linkage disequilibrium (LD) refers to the non-random association of alleles at a pair of genetic loci. It manifests as a deviation of observed haplotype frequencies from the frequencies expected under the assumption that alleles at the 2 loci associate independently. As a fundamental concept in population genetics, LD appears in a wide variety of contexts, such as association mapping and detection of natural selection [1–5].

The original measure of LD for a pair of biallelic loci, one with alleles A and a and the other with alleles B and b , was $D = p_{AB} - p_A p_B$, where p_A and p_B represent the frequencies of alleles A and B , respectively, and p_{AB} is the frequency of the 2-locus haplotype containing alleles A and B [6]. The frequencies p_A and p_B can be measured for the 2 loci separately, each in the absence of information on the other locus, whereas the evaluation of the frequency p_{AB} uses information on the co-occurrence of alleles within individuals at the 2 loci.

Because LD is a property of a relationship between a pair of loci, the values of the allele frequencies at the 2 loci under consideration can affect the potential strength of that relationship. This dependence is a recognized feature

of LD measurement: soon after the initial development of the measure D , the quantity $|D'|$ was introduced as a normalization of D that has the same maximal value irrespective of the allele frequencies at the constituent loci [7].

Many measures of LD have been proposed, each with different arguments favoring its use [1, 3, 8–12]. For example, the popular measure r^2 [13] has the property that it can be interpreted as a squared correlation coefficient between indicator variables for the presence of allele A at the first locus and allele B at the second locus. Each allelic indicator variable is a Bernoulli trial, so that the squared covariance in the numerator of r^2 , D^2 , is obtained by examining the probability that both indicator variables simultaneously equal 1. Features of r^2 in a population evolving according to a standard neutral model are closely related to the population recombination rate $4N_e c$, where N_e is the effective population size and c is the recombination rate between 2 loci [14, 15]. In addition, a calculation of the sample size necessary to detect disease association at a marker locus in LD with a disease locus relies specifically on a measurement of r^2 between the marker and disease loci [3, 16].

The measure d [17] contains an asymmetry between the pair of loci that can be useful if ascertainment of haplotypes forces specific frequencies for one of the loci. This asymmetry is potentially of use in association mapping in the context of a case-control study, where the B/b locus is taken to contain the disease allele, with A/a being a marker locus [9, 18]. In this context, d can also be interpreted as the difference in the proportions of disease and normal alleles found on the same haplotype with a particular marker allele [9].

The measure ρ has been argued to be informative in a model-based perspective on LD, in which it is treated as the probability that a haplotype chosen at random descends without recombination from a population of haplotypes that excludes the aB haplotype [19]. Specifically, given a set of allele and haplotype frequencies, ρ satisfies

$$\rho \begin{bmatrix} p_B & p_A - p_B \\ 0 & 1 - p_A \end{bmatrix} + (1 - \rho) \begin{bmatrix} p_A p_B & p_A(1 - p_B) \\ (1 - p_A)p_B & (1 - p_A)(1 - p_B) \end{bmatrix} = \begin{bmatrix} p_{AB} & p_{Ab} \\ p_{aB} & p_{ab} \end{bmatrix}. \quad (1)$$

A fifth measure, a normalization of r^2 termed r^2/r^2_{\max} [20], has the same property as $|D'|$ in that its maximum is invariant with respect to the values of p_A and p_B .

All of these normalized measures – $|D'|$, r^2 , d , ρ , and r^2/r^2_{\max} – have numerators that are functions of D and denominators that are functions of the single-locus quantities p_A and p_B . The normalizations introduce different consequences for the maximal values of the statistics as functions of p_A and p_B [8, 20, 21]. They also affect the symmetries of

the statistics both with respect to exchanges of the 2 loci and with respect to exchanges of the alleles at one or both loci.

Many applications of LD statistics implement numerical cutoffs to assess if a desired degree of association has been met by a pair of loci, with only those locus pairs whose LD value exceeds the threshold regarded as having done so. For example, pairwise LD thresholds have been used in defining the boundaries of haplotype blocks [22, 23]. They have also been applied to select tag SNP sets to assay in association studies, choosing tags by the number of non-tags with which they achieve a minimum LD cutoff and evaluating the fraction of non-tags that achieve an LD cutoff with at least one tag [24, 25]. LD thresholds have also been employed for such purposes as visualizing tiered LD levels [26], pruning correlated markers in polygenic risk score calculations [27], and generating networks whose vertices represent loci and whose edges connect locus pairs with LD values exceeding a cutoff [28].

The frequent use of pairwise LD thresholds motivates studies of the implicit properties of allele frequencies forced by the thresholds, and more generally, of the way in which numerical values and interpretations of the various statistics depend on allele frequencies. This paper examines such properties and other mathematical features of the various D -based statistics. Although the statistics all range from 0 to 1, owing to their different normalizations and constraints, the meaning of a numerical value of one statistic can differ from the meaning of the same value of another statistic. Our goal is to characterize properties of the range, dependencies, and typical magnitudes of the measures, in order to assist in giving insight about values observed in empirical and theoretical studies of LD.

VanLiere and Rosenberg [20] studied the maximal value of r^2 as a function of p_A and p_B (see also Eberle et al. [29]), in addition to considering such quantities as the mean maximal value of r^2 over the unit square for the pair of allele frequencies, the mean maximal value of r^2 over values of p_B for fixed values of p_A , and the set of permissible values of p_B given r^2 and p_A (see also Wray [30]). With the current emphasis on rare variants in human genetics [31, 32], a salient observation concerning r^2 is that if rare mutations occur at 2 loci on the same common haplotype in different individuals, then r^2 for the pair of loci is likely to have an extremely low value [33, 34], complicating the use of r^2 in comparing LD across locus pairs. Here, we examine aspects of the mathematical properties of LD measures for each of the 5 normalized measures. We also consider the relationships between pairs of measures, finding that some pairs of measures are equal in particular scenarios.

Table 1. The 8 octants in the space of all possible allele frequencies, along with their associated r_{\max}^2 , $|d|_{\max}$, and ρ_{\max} values

Octant	Condition				p_{AB} achieving maximal $ D $	r_{\max}^2	$ d _{\max}$	ρ_{\max}
	$p_A < 1/2$	$p_B < 1/2$	$p_A < p_B$	$p_A + p_B < 1$				
S_1	Yes	No	Yes	No	$p_A + p_B - 1$	$\frac{(1-p_A)(1-p_B)}{p_A p_B}$	$\frac{1-p_A}{p_B}$	-
S_2	No	No	Yes	No	p_A	$\frac{p_A(1-p_B)}{(1-p_A)p_B}$	$\frac{p_A}{p_B}$	-
S_3	No	No	No	No	p_B	$\frac{(1-p_A)p_B}{p_A(1-p_B)}$	-	-
S_4	No	Yes	No	No	$p_A + p_B - 1$	$\frac{(1-p_A)(1-p_B)}{p_A p_B}$	-	1
S_5	No	Yes	No	Yes	0	$\frac{p_A p_B}{(1-p_A)(1-p_B)}$	$\frac{p_A}{1-p_B}$	1
S_6	Yes	Yes	No	Yes	p_B	$\frac{(1-p_A)p_B}{p_A(1-p_B)}$	$\frac{1-p_A}{1-p_B}$	1
S_7	Yes	Yes	Yes	Yes	p_A	$\frac{p_A(1-p_B)}{(1-p_A)p_B}$	-	-
S_8	Yes	No	Yes	Yes	0	$\frac{p_A p_B}{(1-p_A)(1-p_B)}$	-	-

Theory

Setting

We consider 2 biallelic loci, locus 1 with alleles A and a , and locus 2 with alleles B and b . The population frequencies of these alleles are then given by p_A , $p_a = 1 - p_A$, p_B , and $p_b = 1 - p_B$, respectively. Because both p_A and p_B lie in $[0, 1]$, a set of frequencies can be characterized as a point in the unit square with axes p_A and p_B . For ease of notation, following VanLiere and Rosenberg [20], we split this square into octants S_1, S_2, \dots, S_8 , as illustrated in Figure 1. The conditions on p_A and p_B that characterize these octants appear in Table 1. We henceforth assume that the loci are both polymorphic, so that p_A , p_a , p_B , and p_b all lie in $(0, 1)$. The 2 pairs of alleles associate into 4 distinct haplotypes: AB , Ab , aB , and ab , with frequencies p_{AB} , p_{Ab} , p_{aB} , and p_{ab} , respectively (Table 2).

We consider parametric values for the allele frequencies, so that our interest is in LD statistics computed as functions of quantities $p_A, p_a, p_B, p_b, p_{AB}, p_{Ab}, p_{aB},$ and p_{ab} . This setting amounts to considering the statistics in an idealized setting of an infinite population.

The Five Normalized LD Measures

As mentioned earlier, the most basic measure of LD is $D = p_{AB} - p_A p_B$, the difference between the observed frequency of the AB haplotype and its expected frequency under independence of loci 1 and 2. Expressions for D can also be formulated using each of the 3 other possible combinations of alleles at the 2 loci (Ab , aB , and ab). The 4 formulations all give an identical value, up to a change in sign.

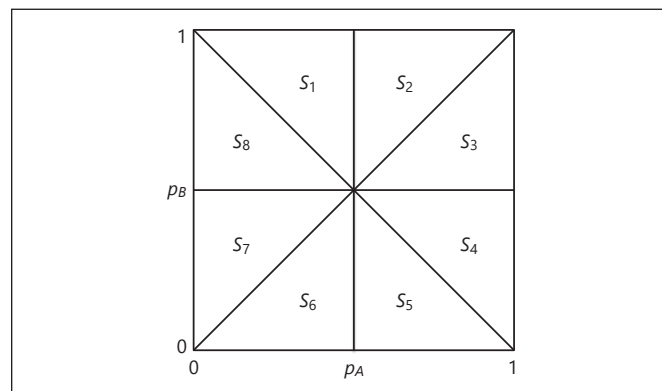


Fig. 1. A unit square showing all possible combinations of the frequencies p_A and p_B . The region is subdivided into 8 octants S_1, \dots, S_8 .

If no association exists between the 2 loci, then we expect $p_{AB} = p_A p_B$, and hence $D = 0$. We consider several LD measures, each of which is a normalization of D . For instance, D' is obtained by normalizing D by its maximal magnitude, given the sign of D :

$$D' = \frac{D}{D_{\max}}$$

where

Table 2. Notation for the allele and haplotype frequencies for a pair of biallelic loci

	Locus 1		Total
	A	a	
Locus 2			
B	p_{AB}	p_{aB}	p_B
b	p_{Ab}	p_{ab}	$1 - p_B$
Total	p_A	$1 - p_A$	1

Table 3. Octants of the allele frequency space in which the different linkage disequilibrium measures can be applied

Octant	D'	r^2	d	ρ	r^2/r^2_{\max}
S_1	Yes	Yes	Yes	No	Yes
S_2	Yes	Yes	Yes	No	Yes
S_3	Yes	Yes	No	No	Yes
S_4	Yes	Yes	No	Yes	Yes
S_5	Yes	Yes	Yes	Yes	Yes
S_6	Yes	Yes	Yes	Yes	Yes
S_7	Yes	Yes	No	No	Yes
S_8	Yes	Yes	No	No	Yes

$$D_{\max} = \begin{cases} \min[p_A(1-p_B), (1-p_A)p_B] & \text{if } D > 0 \\ \min[p_A p_B, (1-p_A)(1-p_B)] & \text{if } D < 0 \end{cases} \quad (2)$$

The r^2 measure is defined as D^2 normalized by the product of all 4 allele frequencies:

$$r^2 = \frac{D^2}{p_A(1-p_A)p_B(1-p_B)} \quad (3)$$

The next 2 LD measures represent the 2 ways in which D can be normalized by the product of 2 of the 4 allele frequencies. If the 2 frequencies represent alleles from the same locus, then we have d , given by Nei and Li [17]:

$$d = \frac{D}{p_B(1-p_B)} \quad (4)$$

By convention, in an association mapping setting, locus 2 is designated as a disease locus, and hence B or b is regarded as a potential disease-causing allele.

If the 2 frequencies instead represent alleles from different loci, then we have ρ , given by Collins and Morton [35]:

$$\rho = \frac{D}{(1-p_A)p_B} \quad (5)$$

Unlike D' and r^2 , both d and ρ introduce an asymmetry in the pair of loci by virtue of the choice of alleles assigned to their denominators.

Lastly, as was noted by VanLiere and Rosenberg [20], the maximal value of r^2 is constrained by the values of the allele frequencies p_A and p_B . Let r^2_{\max} be the maximal value of r^2 possible given p_A and p_B . The measure r^2/r^2_{\max} , introduced by VanLiere and Rosenberg [20], is then simply equal to r^2 normalized by r^2_{\max} .

Prescribed Domains

The measures D' and r^2 can be applied for all values of p_A and p_B in $(0, 1)$, i.e., in all octants in Figure 1. r^2/r^2_{\max} , being derived from r^2 , also has all octants available. However, d and ρ are defined only on part of the domain $(0, 1) \times (0, 1)$. For d , because locus 2 is usually taken to be a disease locus – with one relatively rare allele – and locus 1 is the marker locus, it is assumed that $\min(p_B, p_b) \leq \min(p_A, p_a)$. This assumption restricts d to S_1, S_2, S_5 , and S_6 . For ρ , the allele frequencies are assigned labels such that $D \geq 0$, $p_{aB} \leq p_{Ab}$, $p_B \leq p_A$, and $p_B \leq 1 - p_B$ [19]. Note that $p_{aB} \leq p_{Ab}$ is equivalent to $p_B \leq p_A$, as $p_{aB} = p_B - p_{AB}$ and $p_{Ab} = p_A - p_{AB}$. Together, these conditions restrict the available octants to S_4, S_5 , and S_6 . Domain restrictions are summarized in Table 3, and for d and ρ , we restrict our subsequent analysis to octants in which these measures apply.

Upper Bounds, Mean Maximum Values, and Accessible Regions

We are interested in analyzing the mathematical properties of the 5 LD measures. Because the magnitude of these measures is the quantity of interest, we work with the absolute values $|D'|$ and $|d|$. r^2 and r^2/r^2_{\max} are always non-negative owing to the fact that D^2 is used in their expressions, and ρ is always non-negative because its definition requires $D \geq 0$.

We seek to determine the upper bound, mean maximal value, and accessible region for each of the 5 measures, given the values of p_A and p_B . The mean maximum $\mathbb{E}[m_{\max}]$ of a measure m is defined as its average maximum value, assuming (p_A, p_B) follows a bivariate uniform distribution on the permissible domain over which m applies. Its accessible region for a constant $c \in [0, 1]$, $p_m(c)$, is defined as the proportion of the domain in which the upper bound for the measure is greater than or equal to c .

To determine these mathematical properties, we first must choose a value of p_{AB} that maximizes $|D|$, because all other variables in the expressions for the 5 statistics are fixed given p_A and p_B . The values of p_{AB} that achieve this maximum are the same as those given by VanLiere and Rosenberg [20] for finding the upper bound on r^2 (Fig. 2a), because $|D|$ is maximized if, and only if, D^2 is maximized. Hence, on S_1 and S_4 , the maximum $|D|$ occurs if $p_{AB} = p_A + p_B - 1$; on S_2 and S_7 , if $p_{AB} = p_A$; on S_3 and S_6 , if $p_{AB} = p_B$; and on S_5 and S_8 , if $p_{AB} = 0$. These values appear in Table 1. For all 5 measures, the values of $\mathbb{E}[m_{\max}]$ and $p_m(c)$ appear in Table 4.

$|D'|$: Because $|D'|$ is simply $|D|$ normalized by its maximum value D_{\max} , both its upper bound and the mean maximum $\mathbb{E}[|D'|_{\max}]$ are equal to 1. Furthermore, its accessible region $p_{|D'|}(c)$ is also 1, irrespective of the value of c .

r^2 : The upper bound of r^2 as a function of p_A and p_B , for each octant S_1, \dots, S_8 , was calculated in equations 2–5 of VanLiere and Rosenberg [20]. These results appear in Table 1, and a contour plot of r^2_{\max} is reproduced in Figure 2a. In addition, VanLiere and Rosenberg [20] derived the mean maximum of r^2 , obtaining

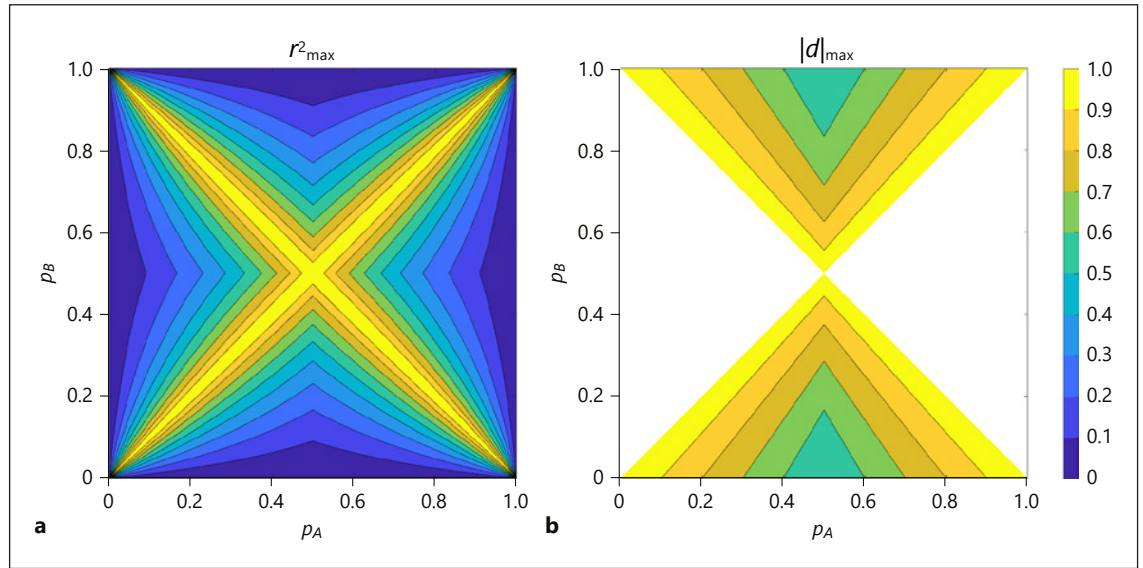


Fig. 2. r^2_{\max} and $|d|_{\max}$ as functions of p_A and p_B . **a** Contour plot of r^2_{\max} . **b** Contour plot of $|d|_{\max}$. The plots consider the maximum over all possible values of p_{AB} . The functions plotted appear in Table 1. $|d|$ is defined only in octants S_1 , S_2 , S_5 , and S_6 .

Table 4. Mean maximum values and accessible regions for the 5 measures

	Mean maximum value	Accessible region
$ D' $	1	1
r^2	$2\pi^2/3 - 4(\ln 2)^2 + 4 \ln 2 - 7 \approx 0.43051$	$1 + \frac{4c}{1-c} + \frac{8c \ln\left(\frac{1}{2} + \frac{1}{2}c\right)}{(1-c)^2}$
$ d $	$3/2 - \ln 2 \approx 0.80685$	1, if $c \leq 0.5$; $\frac{(4c-1)(1-c)}{c}$, if $c > 0.5$
ρ	1	1
r^2/r^2_{\max}	1	1

The mean maximum value of a measure is its average maximum value over its prescribed domain, assuming p_A and p_B are independent and uniformly distributed over the domain.

The accessible region of a measure for a constant $c \in [0, 1]$ is defined as the proportion of the applicable domain in which the upper bound for the measure is greater than or equal to c .

$\mathbb{E}[r^2_{\max}] = 2\pi^2/3 - 4(\ln 2)^2 + 4 \ln 2 - 7 \approx 0.43051$, as well as its accessible region, which is

$$p_{r^2}(c) = 1 + \frac{4c}{1-c} + \frac{8c \ln\left(\frac{1}{2} + \frac{1}{2}c\right)}{(1-c)^2}. \quad (6)$$

A plot of $p_{r^2}(c)$ appears in Figure 3.

$|d|$: By substituting the appropriate value of p_{AB} into the expression for $|d|$, we obtain $|d|_{\max}$ as a function of p_A and p_B in octants S_1 , S_2 , S_5 , and S_6 :

$$S_1 : |d|_{\max}(p_A, p_B) = \frac{|p_A + p_B - 1 - p_A p_B|}{p_B(1-p_B)} = \frac{1-p_A}{p_B} \quad (7)$$

$$S_2 : |d|_{\max}(p_A, p_B) = \frac{p_A - p_A p_B}{p_B(1-p_B)} = \frac{p_A}{p_B} \quad (8)$$

$$S_5 : |d|_{\max}(p_A, p_B) = \frac{|0 - p_A p_B|}{p_B(1-p_B)} = \frac{p_A}{1-p_B} \quad (9)$$

$$S_6 : |d|_{\max}(p_A, p_B) = \frac{p_B - p_A p_B}{p_B(1-p_B)} = \frac{1-p_A}{1-p_B} \quad (10)$$

These results are summarized in Table 1. Figure 2b shows a contour plot of $|d|_{\max}$ in $S_1, S_2, S_5,$ and S_6 , combining equations 7–10. We note some similarities, as well as some differences, with the plot of r^2_{\max} in Figure 2a. Examining the characteristic X shape of the figure, we see that $|d|_{\max}$ can equal 1 if, and only if, the allele frequencies are identical at the 2 loci, $p_A = p_B$ or $p_A = p_b$, as is the case with r^2_{\max} . However, instead of having a symmetric shape over all octants, $|d|_{\max}$ is symmetric with respect to an exchange of p_A and p_a or p_B and p_b , but not with respect to an exchange of p_A and p_B (and thus also p_a and p_b) or p_A and p_b (and thus also p_a and p_B). Its shape is symmetric over $S_1, S_2, S_5,$ and S_6 , the 4 octants on which it can be calculated. Unlike r^2_{\max} , $|d|_{\max}$ does not approach 0 as p_B approaches either 0 or 1. This feature enables $|d|$ to maintain a considerable range of allowable values, even if the minor allele frequency (MAF) at locus 2 is low, as is likely the case in a mapping study in which locus 2 is regarded as causal for a rare disease.

We can quantify the difference in range for $|d|$ and r^2 by comparing the mean maximum value of $|d|$ to that of r^2 . First, we compute volume V_2 , which we define to be the volume of $|d|_{\max}$ over the octant S_2 :

$$\begin{aligned} V_2 &= \int_{\frac{1}{2}}^1 \int_{p_A}^1 \frac{p_A}{p_B} dp_B dp_A \\ &= \frac{3}{16} - \frac{1}{8} \ln 2 \approx 0.10086. \end{aligned} \quad (11)$$

Owing to symmetry, V_2 is equal to corresponding values $V_1, V_5,$ and V_6 . Assuming a uniform joint distribution of p_A and p_B over the octants $S_1, S_2, S_5,$ and S_6 , and noting that these octants have a total area of $1/2$, the mean maximum value of $|d|_{\max}$ is

$$\begin{aligned} \mathbb{E}[|d|_{\max}] &= 4V_2 / \frac{1}{2} \\ &= \frac{3}{2} - \ln 2 \approx 0.80685. \end{aligned} \quad (12)$$

This value exceeds $\mathbb{E}[r^2_{\max}] \approx 0.43051$ derived by VanLiere and Rosenberg [20] under the same assumption of a uniform distribution on the domain, suggesting that $|d|$ can achieve a high magnitude over a considerably larger portion of the allele frequency space than is seen for r^2 . To quantify this difference, we calculate the accessible region $p_{|d|}(c)$. We first focus on S_6 , and then extend the result to the remaining octants using symmetry.

Let A_6 denote the area of the portion of S_6 in which $|d|_{\max} \geq c$. Using eq. 10, the portion of S_6 in which $|d|_{\max} \geq c$ satisfies $p_B \geq (p_A + c - 1)/c$. We now set up an integral to calculate the complement of the desired area, the area of the portion of S_6 in which $|d|_{\max} < c$. Observe from Figure 2b that in S_6 , for $c \geq 1/2$, the horizontal plane $|d|_{\max} = c$ intersects the upper bound at $p_A = 1 - c$ if $p_B = 0$. Therefore, we have

$$\begin{aligned} \frac{1}{8} - A_6 &= \int_{1-c}^1 \int_0^{\frac{p_A+c-1}{c}} 1 dp_B dp_A = \frac{(2c-1)^2}{8c} \\ A_6 &= \frac{(4c-1)(1-c)}{8c}. \end{aligned} \quad (13)$$

For $c \leq 1/2$, all points in octants $S_1, S_2, S_5,$ and S_6 have $|d|_{\max} \geq c$ (Fig. 2b). Hence, in this situation, $p_{|d|}(c)$ is simply 1. For $c \geq 1/2$, applying eq. 13,

$$\begin{aligned} p_{|d|}(c) &= 4A_6 / \frac{1}{2} \\ &= \frac{(4c-1)(1-c)}{c}. \end{aligned} \quad (14)$$

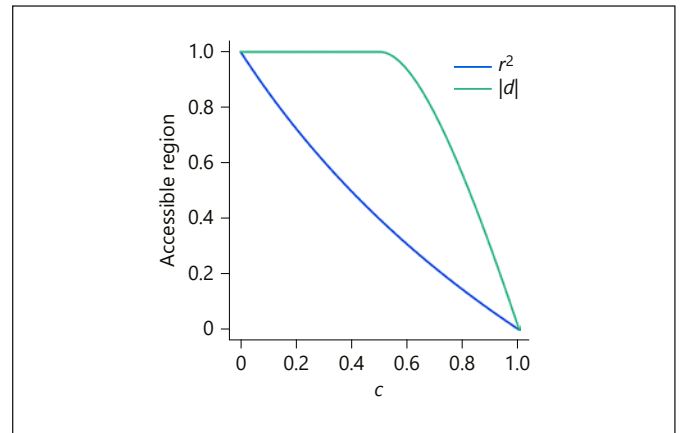


Fig. 3. The portion of the permissible allele frequency space where r^2 and $|d|$ can exceed a specific value of c , as a function of c . $p_{r^2}(c)$ is taken from eq. 6, $p_{|d|}(c)$ is taken from eq. 14, and both functions appear in Table 4.

The piecewise function $p_{|d|}(c)$ appears in Figure 3, alongside a plot of $p_{r^2}(c)$ from eq. 6. The permissible fraction of the frequency space for $|d|$ decreases more slowly as a function of c than does the corresponding function for r^2 .

ρ : For the upper bound on ρ , the following conditions must all be satisfied when assigning labels to the alleles: $D \geq 0, p_B \leq p_A,$ and $p_B \leq 1 - p_B$ [19, 36]. The latter 2 conditions imply that ρ applies only in $S_4, S_5,$ and S_6 . The condition $p_B \leq p_A$ implies $p_B - p_A p_B \leq p_A - p_A p_B$, which in turn implies $(1 - p_A)p_B \leq p_A(1 - p_B)$. This result, in addition to the requirement that $D \geq 0$, indicates that ρ is exactly equal to $|D|$ under the conditions in which ρ applies. Consequently, the upper bound of ρ , its mean maximum $\mathbb{E}[\rho_{\max}]$, and its accessible region $p_\rho(c)$ all equal 1.

r^2/r^2_{\max} : The upper bound, the mean maximum $\mathbb{E}[\{r^2/r^2_{\max}\}_{\max}]$, and the accessible region $p_{r^2/r^2_{\max}}(c)$ all equal 1, by definition of the statistic as r^2 is normalized by r^2_{\max} .

Mean Maximum of r^2 and $|d|$ under a Beta Distribution

In the previous section (“Upper Bounds, Mean Maximum Values, and Accessible Regions”), we examined the mean maximum of the 5 measures, assuming (p_A, p_B) follows a bivariate uniform distribution. For r^2 and $|d|$, the 2 measures that do not have a mean maximum of 1, we can also calculate their mean maximum value under less restrictive assumptions. We now assume p_A and p_B follow independent beta distributions. To preserve symmetry between loci and exchangeability of the alleles at a locus, we consider $p_A, p_B \sim \text{Beta}(\alpha, \alpha)$, and compute $\mathbb{E}[|d|_{\max}]$ and $\mathbb{E}[r^2_{\max}]$ as functions of α .

By analogy with eq. 12, again using octant S_2 , we can set up an integral for $\mathbb{E}[|d|_{\max}]$:

$$\begin{aligned} \mathbb{E}[|d|_{\max}] &= 8 \int_{\frac{1}{2}}^1 \int_{p_A}^1 \frac{p_A}{p_B} \frac{p_A^{\alpha-1}(1-p_A)^{\alpha-1}}{B(\alpha, \alpha)} \frac{p_B^{\alpha-1}(1-p_B)^{\alpha-1}}{B(\alpha, \alpha)} dp_B dp_A \\ &= \frac{8}{[B(\alpha, \alpha)]^2} \int_{\frac{1}{2}}^1 \int_{p_A}^1 \frac{p_A}{p_B} (p_A - p_A^2)^{\alpha-1} (p_B - p_B^2)^{\alpha-1} dp_B dp_A. \end{aligned} \quad (15)$$

Here, $B(\alpha, \alpha) = [\Gamma(\alpha)]^2/\Gamma(2\alpha)$. To compute $\mathbb{E}[r_{\max}^2]$, we use r_{\max}^2 on S_2 (Table 1):

$$\begin{aligned} \mathbb{E}[r_{\max}^2] &= 8 \int_{\frac{1}{2}}^1 \int_{p_A}^1 \frac{p_A(1-p_B) p_A^{\alpha-1} (1-p_A)^{\alpha-1} p_B^{\alpha-1} (1-p_B)^{\alpha-1}}{B(\alpha, \alpha) B(\alpha, \alpha)} dp_B dp_A \\ &= \frac{8}{[B(\alpha, \alpha)]^2} \int_{\frac{1}{2}}^1 \int_{p_A}^1 \frac{p_A(1-p_B)}{(1-p_A)p_B} (p_A - p_A^2)^{\alpha-1} (p_B - p_B^2)^{\alpha-1} dp_B dp_A. \end{aligned} \quad (16)$$

We evaluate eqs. 15 and 16 numerically, using values of α ranging from 0.2 to 5. The results appear in Figure 4. Low values of α imply that the distribution of allele frequencies is skewed toward loci with a low MAF, whereas high α values correspond to greater density in loci with a high MAF. Note that setting $\alpha = 1$ gives a uniform distribution and recovers the values derived in the section ‘‘Upper Bounds, Mean Maximum Values, and Accessible Regions’’ above for the case of $p_A, p_B \sim \text{Uniform}-(0, 1)$.

From Figure 4, we observe that $\mathbb{E}[r_{\max}^2]$ varies considerably as a function of the allele frequency distribution, whereas $\mathbb{E}[|d|_{\max}]$ is more stable as α changes.

The Five Measures as Functions of p_{AB} for Fixed p_A, p_B

For most of our subsequent calculations, to facilitate comparison, we restrict our analysis to values of p_A and p_B in S_6 , as all 5 measures have S_6 in their prescribed domains. Any point not in S_6 can be mapped to a point in S_6 by performing one or more of a set of transformations: (i) a reflection over $p_A = 1/2$ (corresponding to exchanging the p_A and p_a labels), (ii) a reflection over $p_B = 1/2$ (exchanging the p_B and p_b labels), and (iii) a reflection over the $p_A = p_B$ line (exchanging the p_A and p_B labels, and thus also the p_a and p_b labels).

We first compare how each measure varies with the haplotype frequency p_{AB} . Figure 5 illustrates $|D'|$, r^2 , $|d|$, and r^2/r_{\max}^2 as functions of the haplotype frequency p_{AB} , for fixed values of p_A and p_B . In this analysis, ρ is omitted because under the conditions in which it applies, it is exactly equal to $|D'|$. Each of the measures has a value of 0 in the case of linkage equilibrium, at which $p_{AB} = p_A p_B$. Using $p_{AB} = p_A p_B$ as a reference point, we can split the plots for each of the measures into 2 portions: the right arm, where $p_{AB} \geq p_A p_B$ (or $D \geq 0$, corresponding to the case with an excess of haplotypes containing both minor alleles), and the left arm, where $p_{AB} \leq p_A p_B$ (or $D \leq 0$, corresponding to the case with a deficit of haplotypes containing both minor alleles).

$|D'|$: $|D'|$ varies linearly with p_{AB} . However, its left and right arms are in general not symmetric about the line $p_{AB} = p_A p_B$; the absolute value of the derivative of $|D'|$ as a function of p_{AB} differs in the 2 arms. This phenomenon results from the different normalizations applied in obtaining D' , depending on whether D is positive or negative. The left and right arms of $|D'|$ are symmetric only if $p_A = 1/2, p_B = 1/2$, or both. The value of $|D'|$ can always reach 1 irrespective of whether the haplotype containing both minor alleles is in excess or in deficit; in general, no such result holds for the other 3 measures.

r^2 : r^2 varies quadratically as a function of p_{AB} , with the measure increasing at a faster rate the further p_{AB} is from $p_A p_B$. As reported by VanLiere and Rosenberg [20], in S_6 , r^2 can only reach 1 if $p_A = p_B$, and even then, only if an excess of haplotypes containing both minor alleles occurs. Finally, ignoring the truncation imposed by the lower limit of $p_{AB} = 0$, the arms of r^2 are symmetric with respect to the line $p_{AB} = p_A p_B$.

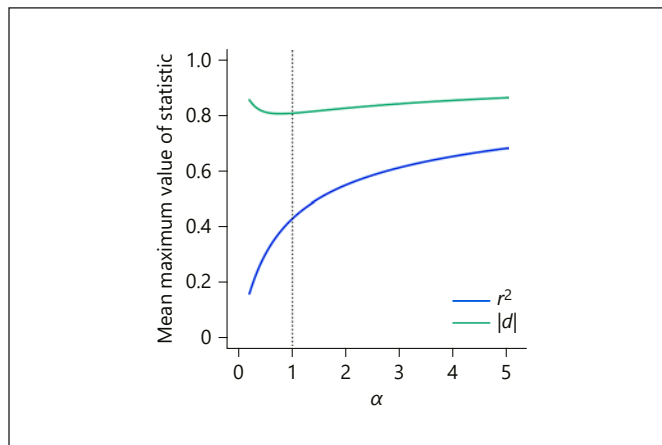


Fig. 4. Mean maximum value of r^2 and $|d|$, if p_A and p_B are drawn from independent Beta- (α, α) distributions. The dotted line indicates $\alpha = 1$, which gives values that are identical to the case in which p_A and p_B are drawn from independent Uniform- $(0, 1)$ distributions.

$|d|$: $|d|$, like $|D'|$, varies linearly as a function of p_{AB} . However, unlike for $|D'|$, the left and right arms of $|d|$ are symmetric, and they have the same absolute value of the derivative as a function of p_{AB} . This pattern occurs because $|d|$ does not necessarily have to reach 1 at the points where p_{AB} lies at its maximum or minimum values, given p_A and p_B . Like r^2 , $|d|$ can only reach 1 on its right arm if $p_A = p_B$. As a result, $|D'| = |d|$ if $p_A = p_B$ and $D \geq 0$, as can be observed from Figure 5.

r^2/r_{\max}^2 : r^2/r_{\max}^2 varies quadratically as a function of p_{AB} , but it increases more quickly compared to r^2 as p_{AB} moves away from $p_A p_B$. It can always reach 1 irrespective of the values of p_A and p_B , but it does so only if the haplotype containing both minor alleles is in excess. If $p_A = p_B$, then $r^2/r_{\max}^2 = r^2$, as the maximum of r^2 is 1 in this case.

Comparison: In general, for all values of p_A, p_B , and p_{AB} , $|D'| \geq |d| \geq r^2$. To demonstrate this, we first show that $|D'| \geq |d|$. Consider $D > 0$, where D is normalized by $\min[p_A(1-p_B), (1-p_A)p_B]$ in the calculation of $|D'|$. If $p_A \geq p_B$, then $p_A(1-p_B) \geq p_B(1-p_B) \geq (1-p_A)p_B$, but if $p_A \leq p_B$, then $p_A(1-p_B) \leq p_B(1-p_B) \leq (1-p_A)p_B$. In either case, $p_B(1-p_B) \geq \min[p_A(1-p_B), (1-p_A)p_B]$, and therefore $|D'| \geq |d|$. A similar calculation shows also that $|D'| \geq |d|$ if $D < 0$.

To see that $|d| \geq r^2$ where $|d|$ applies (S_1, S_2, S_5, S_6), we show $r^2/|d| \leq 1$. We have

$$\frac{r^2}{|d|} = \frac{|D|}{p_A(1-p_A)} = \frac{|p_{AB} - p_A p_B|}{p_A(1-p_A)}. \quad (17)$$

Consider S_6 , where $p_A \geq p_B$ and $p_{AB} = p_B$ maximizes $|D|$ (Table 1):

$$\frac{|p_{AB} - p_A p_B|}{p_A(1-p_A)} \leq \frac{p_B - p_A p_B}{p_A(1-p_A)} = \frac{p_B}{p_A}. \quad (18)$$

Because $p_A \geq p_B$, $r^2/|d| \leq 1$. Eq. 18 and other similar calculations for S_1, S_2 , and S_5 show that $|d| \geq r^2$ for all possible values of p_A and p_B .

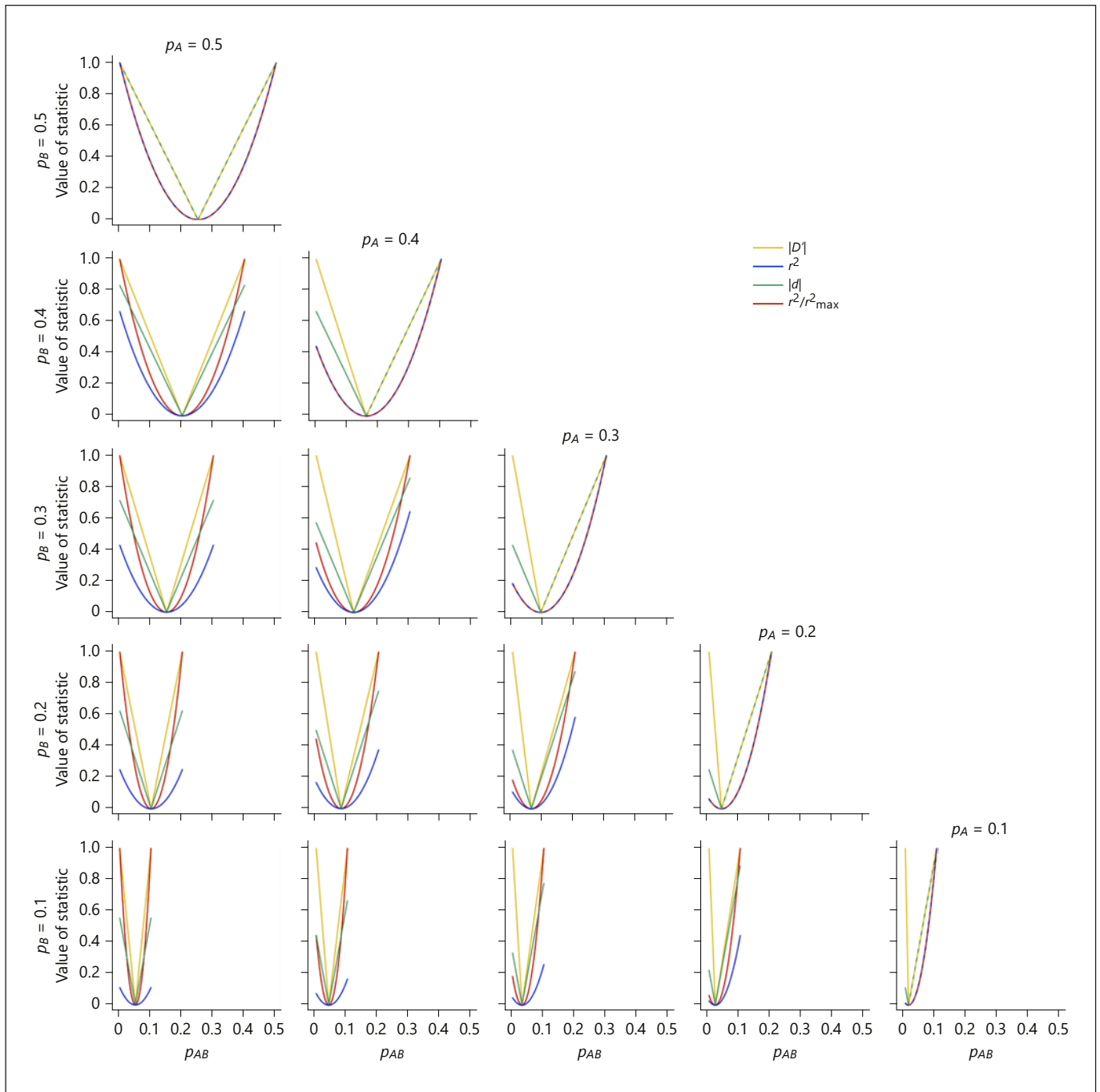


Fig. 5. Values of linkage disequilibrium statistics as functions of the haplotype frequency p_{AB} , for fixed values of p_A and p_B in S_6 .

Mean LD Values Given Fixed p_A and p_B

In the section “Upper Bounds, Mean Maximum Values, and Accessible Regions,” we have described upper bounds of the 5 measures given values of p_A and p_B . We have also computed the mean of the maximum value. Next, we specify a distribution on p_{AB} and calculate the mean values of the measures over the possible

domain. For this analysis, we again use S_6 ; analogous results for other octants follow the same framework. Because $p_B \leq p_A$ in S_6 , p_{AB} lies in $[0, p_B]$. If we further assume $p_{AB} \sim \text{Uniform}(0, p_B)$, then we can compute the mean value of each LD measure as a function of p_A and p_B . For completeness, associated variances are derived in the Appendix.

$|D'|$: $D > 0$ if $p_{AB} > p_A p_B$, and $D < 0$ if $p_{AB} < p_A p_B$. We split the integral for the mean to account for both cases. If p_{AB} is distributed uniformly on $(0, p_B)$, then $\mathbb{E}[|D'|]$ has a constant value of $1/2$ and does not depend on p_A or p_B .

$$\mathbb{E}[|D'|] = \frac{1}{p_B} \left[\int_0^{p_A p_B} \frac{p_A p_B - p_{AB}}{p_A p_B} dp_{AB} + \int_{p_A p_B}^{p_B} \frac{p_{AB} - p_A p_B}{(1-p_A)p_B} dp_{AB} \right] \quad (19)$$

$$= \frac{1}{2}.$$

r^2 : For r^2 , it is not necessary to split the integral.

$$\mathbb{E}[r^2] = \frac{1}{p_B} \int_0^{p_B} \frac{(p_{AB} - p_A p_B)^2}{p_A(1-p_A)p_B(1-p_B)} dp_{AB} \quad (20)$$

$$= \frac{p_B(1-3p_A+3p_A^2)}{3p_A(1-p_A)(1-p_B)}.$$

The result is plotted in Figure 6a.

$|d|$: For d , we again split the integral as we did for $|D'|$.

$$\mathbb{E}[|d|] = \frac{1}{p_B} \left[\int_0^{p_A p_B} \frac{p_A p_B - p_{AB}}{p_B(1-p_B)} dp_{AB} + \int_{p_A p_B}^{p_B} \frac{p_{AB} - p_A p_B}{p_B(1-p_B)} dp_{AB} \right] \quad (21)$$

$$= \frac{1-2p_A+2p_A^2}{2(1-p_B)}.$$

The result is plotted in Figure 6b.

r^2/r^2_{\max} : The result appears in Figure 6c. Note that in octant S_6 , $\mathbb{E}[r^2/r^2_{\max}]$ is a function of only p_A (or in general, allele frequencies at the locus with the higher MAF).

$$\mathbb{E}[r^2/r^2_{\max}] = \frac{1}{p_B} \int_0^{p_B} \frac{\frac{(p_{AB} - p_A p_B)^2}{p_A(1-p_A)p_B(1-p_B)}}{\frac{(1-p_A)p_B}{p_A(1-p_B)}} dp_{AB} \quad (22)$$

$$= \frac{1-3p_A+3p_A^2}{3(1-p_A)^2}.$$

$\mathbb{E}[r^2/r^2_{\max}]$ varies less across the domain for p_A than do $\mathbb{E}[r^2]$ and $\mathbb{E}[|d|]$.

Constraints on One Allele Frequency Given an LD Value and the Allele Frequency at the Other Locus

Here, we examine for each of the 5 LD measures the allowable values of one allele frequency (either p_A or p_B), while fixing the allele frequency at the other locus and specifying the value of the LD measure. Owing to symmetry in the loci, for $|D'|$, r^2 , and r^2/r^2_{\max} , fixing p_A is equivalent to fixing p_B , and we need only examine one case. For $|d|$ and ρ , owing to asymmetries in the formulation of the measures, 2 cases must be considered. For the calculations in this section, we assume – for convenience – that p_A and p_B are the minor allele frequencies (octants S_6 and S_7), and that the constraints on the major allele frequencies will follow accordingly. The results of this section are summarized in Table 5.

Allowable values of p_B , given $|D'|$ and p_A : Having a value of $|D'|$ and a value of p_A does not constrain p_B , as all values of $|D'|$ between 0 and 1 are accessible given a pair of allele frequencies p_A and p_B . This result can be shown by the fact that given p_A and p_B , $|D'|$ is a continuous rational function of p_{AB} . Because 0 and 1 are the extreme values of $|D'|$, by the intermediate value theorem, $|D'|$ can take on any value in $[0, 1]$ (also see Fig. 5).

Allowable values of p_B , given r^2 and p_A : Assuming that p_A , $p_B \leq 1/2$, given p_A and r^2 , the constraint on p_B is

$$\frac{r^2 p_A}{1+r^2 p_A - p_A} \leq p_B \leq \min\left(\frac{1}{2}, \frac{p_A}{r^2 p_A + p_A}\right). \quad (23)$$

This result has been previously reported in eqs. 10 and 11 of VanLiere and Rosenberg [20] and Table 2 of Wray [30].

Allowable values of p_B , given $|d|$ and p_A : Because $|d|$ is not symmetric in the 2 loci, we first assume $|d|$ and p_A are specified, and solve for the range of p_B . Recalling that d applies only in S_1 , S_2 , S_5 , and S_6 , and assuming $p_A, p_B \leq 1/2$, we consider S_6 . From eq. 10,

$$|d| \leq \frac{1-p_A}{1-p_B} \quad (24)$$

$$p_B \geq \frac{p_A + |d| - 1}{|d|}.$$

Taking into account $0 \leq p_B \leq p_A \leq 1/2$, we have

$$\max\left(0, \frac{p_A + |d| - 1}{|d|}\right) \leq p_B \leq p_A. \quad (25)$$

Allowable values of p_A , given $|d|$ and p_B : Next, we assume $|d|$ and p_B are specified, and solve for the range of p_A . In S_6 , from eq. 10,

$$|d| \leq \frac{1-p_A}{1-p_B} \quad (26)$$

$$p_A \leq 1 - |d| + |d| p_B.$$

Taking into account $0 \leq p_B \leq p_A \leq 1/2$, we have

$$p_B \leq p_A \leq \min\left(\frac{1}{2}, 1 - |d| + |d| p_B\right). \quad (27)$$

The upper bound here corresponds to the upper bound for the “frequency difference” measure of LD reported in Table 2 of Wray [30], noting that labels p_A and p_B are reversed in that study. However, a difference exists between the reported lower bounds, which can be attributed to the fact that $p_A \leq 1/2$ is not mandated by Wray [30].

Allowable values of p_B , given ρ and p_A : Because all values of ρ in $[0, 1]$ can be reached with any given set of allele frequencies in the permissible domain, no additional constraint exists on p_B given ρ and p_A .

Allowable values of p_A , given ρ and p_B : For the same reason as in the case in which p_A is instead specified, no additional constraint exists on p_A given ρ and p_B .

Allowable values of p_B , given r^2/r^2_{\max} and p_A : Being given a value of r^2/r^2_{\max} and p_A does not constrain the values p_B can take, as all values of r^2/r^2_{\max} in $[0, 1]$ are accessible for a given set of allele frequencies p_A and p_B .

Data Illustration

We now examine how LD distributions from data, as given by the various measures, relate to our bounds. We use the 1000 Genomes Project (data at <http://csg.sph.umich.edu/abecasis/MACH/download/1000G-PhaseI-Interim.html>), considering LD values on chromosome 22

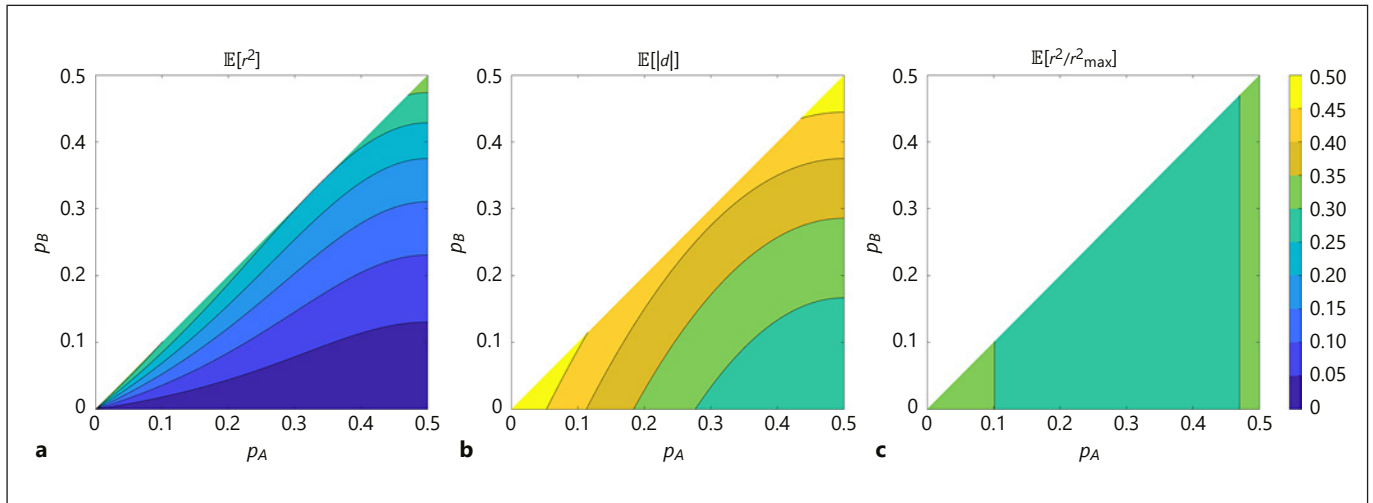


Fig. 6. Mean value of 3 linkage disequilibrium statistics in S_6 as functions of p_A and p_B , assuming $p_{AB} \sim \text{Uniform}(0, p_B)$. **a** r^2 . **b** $|d|$. **c** r^2/r^2_{\max} .

Table 5. Constraints on one allele frequency given an LD value and the allele frequency at the other locus, for the 5 measures

	Fixed allele frequency	
	p_A	p_B
$ D' $	$0 \leq p_B \leq 1$	$0 \leq p_A \leq 1$
r^2	$\frac{r^2 p_A}{1 + r^2 p_A - p_A} \leq p_B \leq \min\left(\frac{1}{2}, \frac{p_A}{r^2 - r^2 p_A + p_A}\right)$	$\frac{r^2 p_B}{1 + r^2 p_B - p_B} \leq p_A \leq \min\left(\frac{1}{2}, \frac{p_B}{r^2 - r^2 p_B + p_B}\right)$
$ d $	$\max\left(0, \frac{p_A + d - 1}{ d }\right) \leq p_B \leq p_A$	$p_B \leq p_A \leq \min\left(\frac{1}{2}, 1 - d + d p_B\right)$
ρ	$0 \leq p_B \leq 1$	$0 \leq p_A \leq 1$
r^2/r^2_{\max}	$0 \leq p_B \leq 1$	$0 \leq p_A \leq 1$

Here, we assume $p_A, p_B \leq 1/2$. Owing to symmetry in the loci, for $|D'|$, r^2 , and r^2/r^2_{\max} , fixing p_A is equivalent to fixing p_B .

of its pooled European population, consisting of 381 individuals: 87 Utah residents of Northern and Western European ancestry, 93 Finnish from Finland, 89 British from England and Scotland, 14 Iberians from Spain, and 98 Toscani from Italy. To ensure inclusion of locus pairs with substantial LD, calculations are restricted to pairs of loci that lie at most 1,000 bp apart. Once again, for ease of comparison, all pairs of allele frequency values outside S_6 are mapped to corresponding points within S_6 .

Pairs of Loci for Which $|D'| = 1$

From 225,159 loci (of 494,975 loci in total) biallelic in the European data, we obtain 1,742,020 pairs of loci separated by at most 1,000 bp. Of these, 1,465,140 pairs have

$|D'| = 1$, indicating the presence of only 2 or 3 of the 4 possible haplotypes. Recalling that $|D'|$ can reach 1 either if an excess or a deficit of haplotypes containing both minor alleles occurs, 349,837 locus pairs belong to the former case, and 1,115,303 to the latter.

r^2 : We now examine on S_6 how r^2 is distributed in relation to the upper bound. If $|D'| = 1$ for a pair of loci, then the r^2 value lies on one of 2 surfaces. If $|D'| = 1$ as the result of an excess of haplotypes containing both minor alleles, corresponding to $p_{AB} = p_B$, then r^2 lies on the surface that defines the upper bound on S_6 , or

$$r^2 = \frac{(1 - p_A)p_B}{p_A(1 - p_B)} \quad (28)$$

as given by eq. 4 in VanLiere and Rosenberg [20]. In S_6 , with a deficit of haplotypes containing both minor alleles, $|D'| = 1$ is achieved if $p_{AB} = 0$, which results in the r^2 surface

$$r^2 = \frac{(0 - p_A p_B)^2}{p_A(1 - p_A)p_B(1 - p_B)} = \frac{p_A p_B}{(1 - p_A)(1 - p_B)}. \quad (29)$$

The surfaces and data points for these 2 cases appear in Figure 7a and b.

|d|: As was seen with r^2 , $|d|$ for a locus pair lies on one of 2 surfaces if $|D'| = 1$ (Fig. 7c, d). Once again, if $p_{AB} = p_B$, then the associated surface is the upper bound of $|d|$ on S_6 , as in eq. 10. If $p_{AB} = 0$, then the points lie on

$$|d| = \frac{|0 - p_A p_B|}{p_B(1 - p_B)} = \frac{p_A}{1 - p_B}. \quad (30)$$

r^2/r_{\max}^2 : Finally, we repeat the analysis for r^2/r_{\max}^2 values. If $p_{AB} = p_B$, then $r^2 = r_{\max}^2$, and therefore $r^2/r_{\max}^2 = 1$. If instead $p_{AB} = 0$, then using eqs. 28 and 29, we obtain

$$\frac{r^2}{r_{\max}^2} = \frac{\frac{p_A p_B}{(1 - p_A)(1 - p_B)}}{\frac{p_A(1 - p_B)}{p_A(1 - p_B)}} = \left(\frac{p_A}{1 - p_A}\right)^2. \quad (31)$$

The surfaces and points corresponding to these 2 cases appear in Figure 7e and f.

Pairs of Loci for Which $|D'| < 1$

In the special case in which $|D'| = 1$ for a pair of loci, we have seen that corresponding values for r^2 , $|d|$, and r^2/r_{\max}^2 lie on well-defined surfaces. Although most locus pairs from our data fall within this category, for 276,880 of 1,742,020 pairs, $|D'| < 1$. For these pairs, we examine how the values of $|D'|$, r^2 , $|d|$, and r^2/r_{\max}^2 are distributed within their ranges.

Recognizing that the 4 measures can exhibit different distribution patterns at different allele frequencies, we sample pairs of loci for which p_A and p_B have MAF values within 4 specified ranges, representing very low, low, intermediate, and high MAF: (0, 0.02], [0.04, 0.06], [0.24, 0.26], and [0.44, 0.46]. Distributions of the values for the measures appear as a series of histograms in Figure 8, with each panel representing one pair of allele frequency ranges for p_A and p_B ; because $p_B \leq p_A$ in S_6 , only 10 of the 16 possible combinations of joint allele frequency ranges are possible.

From Figure 8, we can observe a few properties of the distributions. First, in accordance with our theoretical results, the range of values for r^2 and $|d|$ does not extend to 1 in the panels that are off the diagonal, where $p_A = p_B$ is not possible. In particular, the limitation on the range of

r^2 is more pronounced than that of $|d|$, when comparing within similar allele frequency ranges. This constraint also results in a large number of r^2 values being close to 0, especially if p_B is small (bottom row).

In addition, although we selected only locus pairs which are at most 1,000 bp apart, if p_B is small, then relatively few pairs have a high LD value. This result holds especially for r^2 , but also for measures such as $|D'|$ and r^2/r_{\max}^2 that always have upper bound 1. If one of the loci has a low MAF, then small changes in the haplotype frequency can have large effects on LD measures, especially those normalized to potentially reach 1 irrespective of the marginal allele frequencies (see also the $p_B = 0.1$ panels of Fig. 5).

Comparing scenarios on the diagonal, where it is possible in principle for all 4 measures to achieve high LD values, we see that high LD values are more frequently observed for high and intermediate MAF than for low and very low MAF. For pairs of loci with low MAF, it is unusual for haplotypes to contain the rare allele at both loci, as the rare alleles likely result from relatively recent mutations that have taken place on the same common haplotype, but in different individuals. Thus, because the rare alleles are unlikely to co-occur, the nature of evolutionary descent makes it improbable that the LD-maximizing scenario that couples the rare variants will be obtained. These considerations support a cautious perspective when interpreting LD measures in the case that one or both loci have a low MAF.

Discussion

In this paper, we have described the domains of 5 LD measures that are defined by normalizations of D or D^2 , with a function of p_A and p_B in the denominator. Based on these domains, we have calculated the upper bound, mean maximum, and accessible region for each of the 5 measures. Three of the measures ($|D'|$, ρ , and r^2/r_{\max}^2) can be considered “unrestricted,” in that their upper bound, mean maximum, and accessible region are all equal to 1. However, for the remaining 2 measures (r^2 and $|d|$), these values depend on the allele frequencies of the pair of loci under consideration.

For each of the 5 measures, its description, proposed usages, and mathematical properties are summarized in Table 6. The table provides examples illustrating how a measure’s mathematical properties can inform its use. For instance, $|d|$ allows for a theoretically wider range of values compared to r^2 , with a mean maximum of 0.80685 compared to 0.43051 for r^2 . The increased range of $|d|$ is evident in the analysis of genetic data, which suggests that

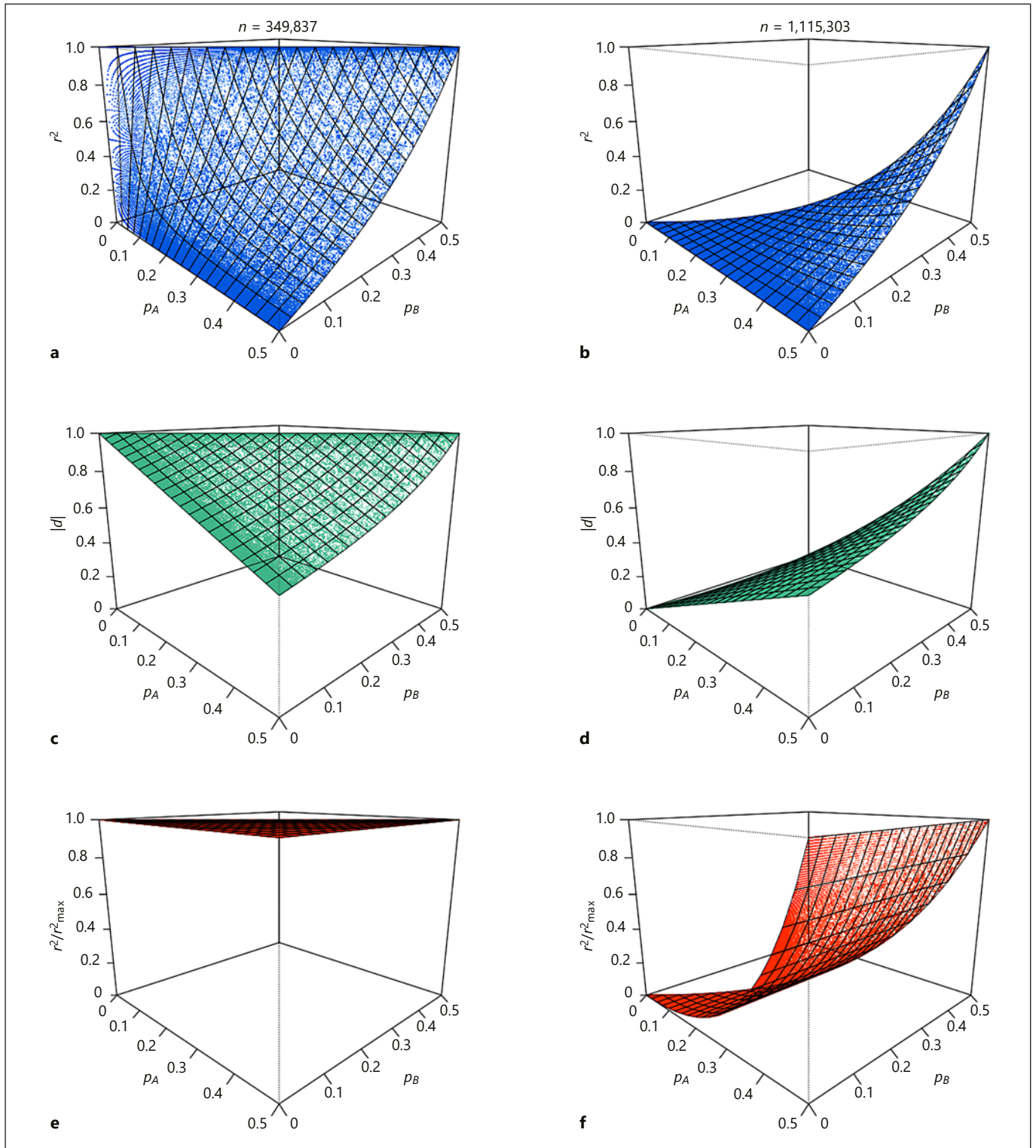


Fig. 7. The distributions of r^2 , $|d|$, and r^2/r^2_{\max} values calculated from data in S_6 , when $|D'| = 1$. **a** r^2 values, if $p_{AB} = p_B$, lying on the surface $r^2 = (1 - p_A)p_B/[p_A(1 - p_B)]$. **b** r^2 values, if $p_{AB} = 0$, lying on the surface $r^2 = p_A p_B / [(1 - p_A)(1 - p_B)]$. **c** $|d|$ values, if $p_{AB} = p_B$, lying on the surface $|d| = (1 - p_A)/(1 - p_B)$. **d** $|d|$ values, if $p_{AB} = 0$, lying on the surface $|d| = p_A/(1 - p_B)$. **e** r^2/r^2_{\max} values, if $p_{AB} = p_B$, lying on the surface $r^2/r^2_{\max} = 1$. **f** r^2/r^2_{\max} values, if $p_{AB} = 0$, lying on the surface $r^2/r^2_{\max} = [p_A/(1 - p_A)]^2$.

Fig. 8. The distributions of values of linkage disequilibrium statistics calculated from data in S_6 , when $|D'| < 1$, given specific ranges of values for p_A and p_B . For each of 4 windows for p_A and 4 windows for p_B , we plot points in histograms based on their values of each of 4 statistics ($|D'|$, r^2 , $|d|$, and r^2/r_{\max}^2). The number of locus pairs falling into a pair of windows appears in the top right corner of the group of 4 histograms associated with the window pair.

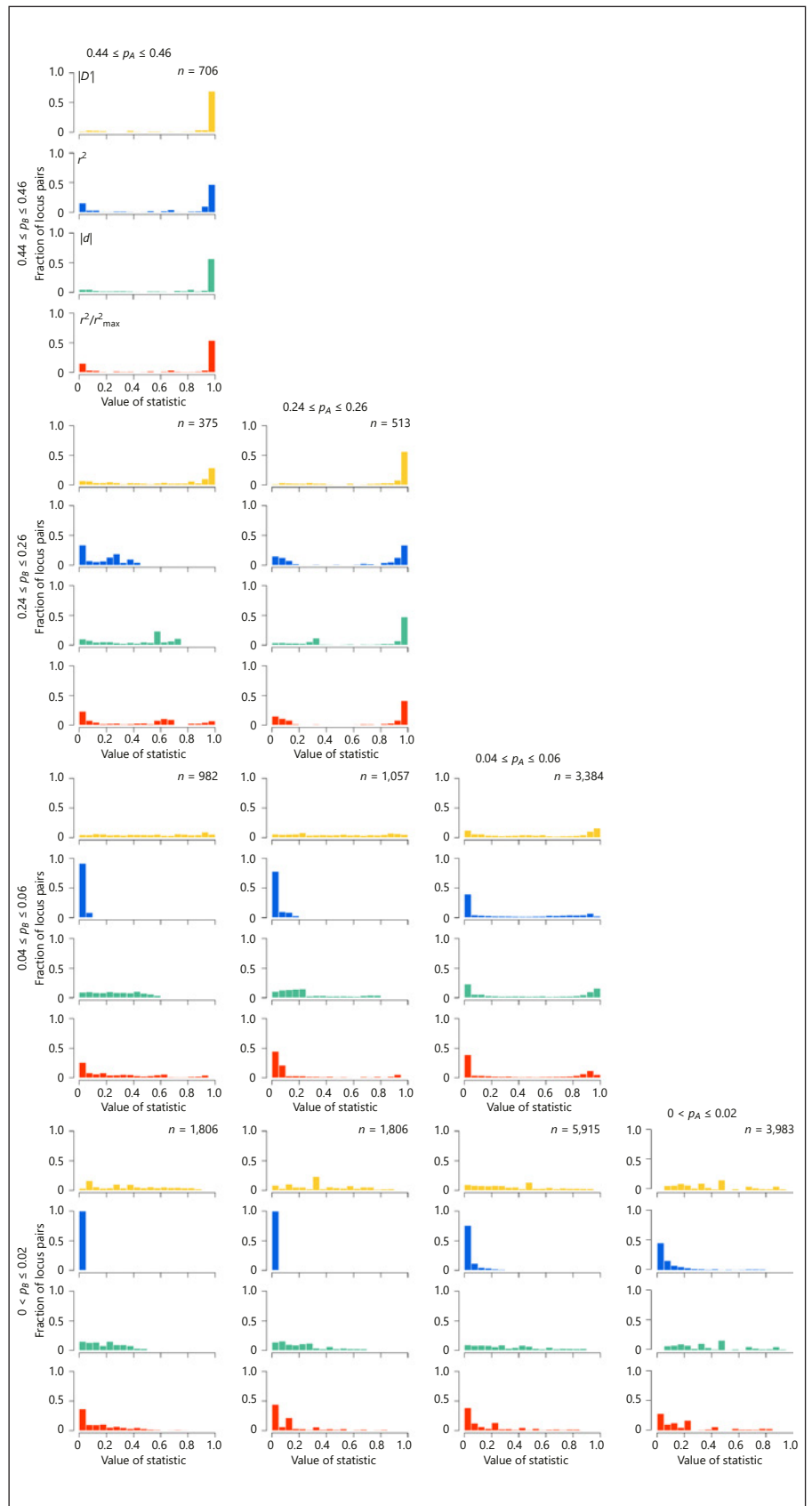


Table 6. Description, usages, and mathematical properties of the 5 LD measures for biallelic loci

Statistic	Description	Noted usages in the literature	Mathematical properties
$ D' $	Normalization of $ D $ by its theoretical maximum value for a given set of allele frequencies	Detecting “complete” LD (where one of the four haplotypes is absent), an indication of whether recombination has occurred between the two loci [11]	$ D' $ varies linearly as a function of p_{AB} (Fig. 5) Upper bound of 1 for all allele frequencies Assuming a uniform distribution of p_{AB} over the range of values it can take, the mean and variance of $ D' $ are both constant values (eqs. 19 and 33)
r^2	Squared correlation coefficient measure between allelic indicator variables	Testing for independence between a pair of loci by a χ^2 test [16] Association studies, where a mathematical relationship exists between r^2 and the sample size needed to detect association between a marker and disease phenotype [16]	r^2 varies quadratically as a function of p_{AB} (Fig. 5) Low upper bound and small range of values if MAF is low (Fig. 2a) Mean maximum value varies considerably as a function of the allele frequency distribution (Fig. 4)
$ d $	Difference in the proportions of disease and normal alleles found on the same haplotype with a particular marker allele	Association mapping for rare diseases in which case-control sampling is employed [9]	$ d $ varies linearly as a function of p_{AB} (Fig. 5) Upper bound has an intermediate value; measure still has a considerable range even at low MAF (Fig. 2b) Mean maximum value relatively stable as a function of the allele frequency distribution (Fig. 4)
ρ	Probability that a haplotype chosen at random descends without recombination from a population of haplotypes that excludes one of the four possible haplotypes	Mapping of marker association and localization of disease loci [19]	Identical to $ D' $ in the octants in which it can be applied
r^2/r^2_{\max}	Normalization of r^2 by its theoretical maximum value for a given set of allele frequencies	If a range that is independent of allele frequency is desired, but the measure still maintains some connection to r^2 [20]	r^2/r^2_{\max} varies quadratically as a function of p_{AB} (Fig. 5) Upper bound of 1 for all allele frequencies Assuming a uniform distribution of p_{AB} over the range of values it can take, the mean and variance of r^2/r^2_{\max} both depend only on the locus with the larger MAF (eqs. 22 and 39)

LD, linkage disequilibrium; MAF, minor allele frequency.

empirical $|d|$ values are more differentiated than corresponding r^2 values (Fig. 8).

In a sense, $|d|$ can be considered a measure that is “intermediate” between $|D'|$ and r^2 . First, its value always lies between r^2 and $|D'|$. It also has properties in common each with r^2 and $|D'|$. Like $|D'|$, given p_A and p_B , $|d|$ varies linearly as a function of p_{AB} , possessing a property that r^2 does not share. However, like r^2 but unlike $|D'|$, $|d|$ is symmetric in p_{AB} around the linkage equilibrium value $p_{AB} = p_A p_B$ (Fig. 5).

We have also identified situations in which some of these measures are equal to one another. Among the measures, ρ uniquely requires $D \geq 0$. This requirement, along with the conditions $p_B \leq p_A$ and $p_B \leq 1 - p_B$, can be satisfied by (i) a reflection over $p_A = 1/2$ (exchanging the p_A and p_a labels), (ii) a reflection over $p_B = 1/2$ (exchanging p_B and p_b), or (iii) both. Under these prescribed conditions for the use of ρ , it exactly equals $|D'|$, a fact that had also been noted by Shete [36] and Mangin et al. [37]. Furthermore, $|d|$ equals $|D'|$ if $p_A = p_B$ and the haplotype

containing both minor alleles is in excess. This result can be seen by observing the right arms of $|D'|$ and $|d|$ in plots along the diagonal of Figure 5, and also from noting that in the right arm, where $D \geq 0$, the normalizations $\min[p_A(1-p_B), (1-p_A)p_B]$ and $p_B(1-p_B)$ for $|D'|$ and $|d|$, respectively, agree if $p_A = p_B$.

By quantifying the degree to which values for the different LD statistics change in response to shifts in allele and haplotype frequencies, the results provide context to the use of LD thresholds in various statistical genetics applications. Some uses impose thresholds in LD measures alongside minimal-MAF cutoffs [24, 25, 27], and our results can be used to understand the behavior of the statistics in permissible ranges specified by simultaneous LD and MAF thresholds. Additionally, the results are useful for low-MAF loci, for which the rare alleles are unlikely to occur on the same haplotype. In particular, they illustrate that r^2 is the most tightly constrained measure (Fig. 5, eq. 18), so that other measures might provide a broader range of values when computing LD statistics for loci with rare variants.

We note that we have focused on parametric aspects of LD measures rather than LD estimated from samples. Sampling properties can be examined, both in models that view alleles as draws from a parametric allele frequency distribution and in coalescent perspectives whose allele frequencies represent outcomes of a generative model (e.g., Weir [38], Rosenberg and Blum [39], Song and Song [40]). The functional forms of estimators can then potentially be combined with bounds on parametric LD measures to produce corresponding bounds on the estimators (e.g., Alcalá and Rosenberg [41], p. 1590).

Many other measures of LD exist that are not included in our analysis [1, 3, 8–11]. Other measures are sometimes normalized by a quantity that includes a haplotype frequency, rather than a function of allele frequencies only, and thus do not lend themselves well to the framework in this paper. Similarly, LD measures used specifically in cases pertaining to multiallelic loci, such as the multiallelic $|D'|$ [8, 42, 43], require additional parameters. Of the LD measures that are described by a ratio of a function of D to a product of allele frequencies for biallelic loci, we have taken a comprehensive look at the most natural statistics with that form.

We initially assumed Uniform-(0, 1) distributions on the allele frequencies to perform computations for the mean maximum of the various measures. This choice, as in VanLiere and Rosenberg [20], permits us to obtain mathematical insight into those measures across their prescribed ranges. In some applications, weighted distributions, such as the Beta-(α, α) distribution we subsequently used, can be applied in place of the uniform distribution.

Acknowledgments

We thank M. Edge for discussions.

Statement of Ethics

The authors have no ethical conflicts to disclose.

Disclosure Statement

The authors have no conflicts of interest to declare.

Funding Sources

We acknowledge NIH grant R01 HG005855 for support.

Author Contributions

J.T.L.K. and N.A.R. conceived the study, performed the mathematical computations, analyzed the results, and wrote the manuscript. J.T.L.K. performed the data analysis.

Appendix

Here, we provide the variances of $|D'|$, r^2 , $|d|$, and r^2/r^2_{\max} , under the assumptions in the section “Mean LD Values Given Fixed p_A and p_B ” above. For computing the variances, we use eqs. 19–22 to supply the means $\mathbb{E}[|D'|]$, $\mathbb{E}[r^2]$, $\mathbb{E}[|d|]$, and $\mathbb{E}[r^2/r^2_{\max}]$, respectively.

$$\mathbb{E}[|D'|^2] = \frac{1}{p_B} \left[\int_0^{p_A p_B} \frac{(p_A p_B - p_{AB})^2}{p_A^2 p_B^2} dp_{AB} + \int_{p_A p_B}^{p_B} \frac{(p_{AB} - p_A p_B)^2}{(1-p_A)^2 p_B^2} dp_{AB} \right] = \frac{1}{3}. \quad (32)$$

$$\text{Var}[|D'|] = \frac{1}{3} - \left(\frac{1}{2}\right)^2 = \frac{1}{12}. \quad (33)$$

$$\mathbb{E}[r^4] = \frac{1}{p_B} \int_0^{p_B} \frac{(p_{AB} - p_A p_B)^4}{p_A^2 (1-p_A)^2 p_B^2 (1-p_B)^2} dp_{AB} = \frac{[p_A^5 + (1-p_A)^5] p_B^2}{5 p_A^2 (1-p_A)^2 (1-p_B)^2}. \quad (34)$$

$$\text{Var}[r^2] = \frac{[p_A^5 + (1-p_A)^5] p_B^2}{5 p_A^2 (1-p_A)^2 (1-p_B)^2} - \left[\frac{(1-3p_A+3p_A^2)p_B}{3p_A(1-p_A)(1-p_B)} \right]^2 = \frac{(4-15p_A+15p_A^2)p_B^2}{45p_A^2(1-p_A)^2(1-p_B)^2}. \quad (35)$$

$$\mathbb{E}[|d|^2] = \frac{1}{p_B} \int_0^{p_B} \frac{(p_{AB} - p_A p_B)^2}{p_B^2 (1-p_B)^2} dp_{AB} \quad (36)$$

$$= \frac{1-3p_A+3p_A^2}{3(1-p_B)^2}$$

$$\text{Var}[|d|] = \frac{1-3p_A+3p_A^2}{3(1-p_B)^2} - \left[\frac{1-2p_A+2p_A^2}{2(1-p_B)} \right]^2 \quad (37)$$

$$= \frac{1-12p_A^2+24p_A^3-12p_A^4}{12(1-p_B)^2}$$

$$\mathbb{E}\left[\left(r^2 / r_{\max}^2\right)^2\right] = \frac{1}{p_B} \int_0^{p_B} \left[\frac{\frac{(p_{AB} - p_A p_B)^2}{p_A(1-p_A)p_B(1-p_B)}}{\frac{(1-p_A)p_B}{p_A(1-p_B)}} \right]^2 dp_{AB} \quad (38)$$

$$= \frac{p_A^5 + (1-p_A)^5}{5(1-p_A)^4}$$

$$\text{Var}\left[r^2 / r_{\max}^2\right] = \frac{p_A^5 + (1-p_A)^5}{5(1-p_A)^4} - \left[\frac{1-3p_A+3p_A^2}{3(1-p_A)^2} \right]^2 \quad (39)$$

$$= \frac{4-15p_A+15p_A^2}{45(1-p_A)^4}$$

References

- Hudson RR. Linkage disequilibrium and recombination. In: Balding DJ, Bishop M, Cannings C, editors. *Handbook of Statistical Genetics*. Chichester: Wiley; 2001. pp. 309–24.
- Nordborg M, Tavaré S. Linkage disequilibrium: what history has to tell us. *Trends Genet*. 2002 Feb;18(2):83–90.
- McVean G. Linkage disequilibrium, recombination and selection. In: Balding DJ, Bishop M, Cannings C, editors. *Handbook of Statistical Genetics*. 3rd ed. Chichester: Wiley; 2007. pp. 909–44.
- Slatkin M. Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nat Rev Genet*. 2008 Jun; 9(6):477–85.
- Weir BS. Linkage disequilibrium and association mapping. *Annu Rev Genomics Hum Genet*. 2008;9(1):129–42.
- Lewontin RC, Kojima K. The evolutionary dynamics of complex polymorphisms. *Evolution*. 1960;14:458–72.
- Lewontin RC. The interaction of selection and linkage. I. General considerations; heterotic models. *Genetics*. 1964 Jan;49(1):49–67.
- Hedrick PW. Gametic disequilibrium measures: proceed with caution. *Genetics*. 1987 Oct;117(2):331–41.
- Devlin B, Risch N. A comparison of linkage disequilibrium measures for fine-scale mapping. *Genomics*. 1995 Sep;29(2):311–22.
- Sabatti C, Risch N. Homozygosity and linkage disequilibrium. *Genetics*. 2002 Apr;160(4): 1707–19.
- Mueller JC. Linkage disequilibrium for different scales and applications. *Brief Bioinform*. 2004 Dec;5(4):355–64.
- Zapata C. On the uses and applications of the most commonly used measures of linkage disequilibrium from the comparative analysis of their statistical properties. *Hum Hered*. 2011;71(3):186–95.
- Hill WG, Robertson A. Linkage disequilibrium in finite populations. *Theor Appl Genet*. 1968 Jun;38(6):226–31.
- Hudson RR. The sampling distribution of linkage disequilibrium under an infinite allele model without selection. *Genetics*. 1985 Mar; 109(3):611–31.
- McVean GA. A genealogical interpretation of linkage disequilibrium. *Genetics*. 2002 Oct; 162(2):987–91.
- Pritchard JK, Przeworski M. Linkage disequilibrium in humans: models and data. *Am J Hum Genet*. 2001 Jul;69(1):1–14.
- Nei M, Li WH. Non-random association between electromorphs and inversion chromosomes in finite populations. *Genet Res*. 1980 Feb;35(1):65–83.
- Kaplan N, Weir BS. Expected behavior of conditional linkage disequilibrium. *Am J Hum Genet*. 1992 Aug;51(2):333–43.
- Morton NE, Zhang W, Taillon-Miller P, Ennis S, Kwok PY, Collins A. The optimal measure of allelic association. *Proc Natl Acad Sci USA*. 2001 Apr;98(9):5217–21.
- VanLiere JM, Rosenberg NA. Mathematical properties of the r^2 measure of linkage disequilibrium. *Theor Popul Biol*. 2008 Aug; 74(1):130–7.
- Lewontin RC. On measures of gametic disequilibrium. *Genetics*. 1988 Nov;120(3):849–52.
- Gabriel SB, Schaffner SF, Nguyen H, Moore JM, Roy J, Blumenstiel B, et al. The structure of haplotype blocks in the human genome. *Science*. 2002 Jun;296(5576):2225–9.
- Wall JD, Pritchard JK. Haplotype blocks and linkage disequilibrium in the human genome. *Nat Rev Genet*. 2003 Aug;4(8):587–97.
- Carlson CS, Eberle MA, Rieder MJ, Yi Q, Kruglyak L, Nickerson DA. Selecting a maximally informative set of single-nucleotide polymorphisms for association analyses using linkage disequilibrium. *Am J Hum Genet*. 2004 Jan;74(1):106–20.
- de Bakker PI, Yelensky R, Pe'er I, Gabriel SB, Daly MJ, Altshuler D. Efficiency and power in genetic association studies. *Nat Genet*. 2005 Nov;37(11):1217–23.
- Barrett JC, Fry B, Maller J, Daly MJ. Haplotype: analysis and visualization of LD and haplotype maps. *Bioinformatics*. 2005 Jan; 21(2):263–5.
- Vilhjálmsdóttir BJ, Yang J, Finucane HK, Gussev A, Lindström S, Ripke S, et al.; Schizophrenia Working Group of the Psychiatric Genomics Consortium, Discovery, Biology, and Risk of Inherited Variants in Breast Cancer (DRIVE) study. Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *Am J Hum Genet*. 2015 Oct;97(4): 576–92.
- Kempainen P, Knight CG, Sarma DK, Hlaing T, Prakash A, Maung Maung YN, et al. Linkage disequilibrium network analysis (LDna) gives a global view of chromosomal inversions, local adaptation and geographic structure. *Mol Ecol Resour*. 2015 Sep;15(5): 1031–45.
- Eberle MA, Rieder MJ, Kruglyak L, Nickerson DA. Allele frequency matching between SNPs reveals an excess of linkage disequilibrium in genic regions of the human genome. *PLoS Genet*. 2006 Sep;2(9):e142.
- Wray NR. Allele frequencies and the r^2 measure of linkage disequilibrium: impact on design and interpretation of association studies. *Twin Res Hum Genet*. 2005 Apr;8(2):87–94.
- Auer PL, Lettre G. Rare variant association studies: considerations, challenges and opportunities. *Genome Med*. 2015 Feb;7(1):16.
- Bomba L, Walter K, Soranzo N. The impact of rare and low-frequency genetic variants in common disease. *Genome Biol*. 2017 Apr; 18(1):77.
- Li B, Liu DJ, Leal SM. Identifying rare variants associated with complex traits via sequencing. *Curr Protoc Hum Genet*. 2013 Jul;78:1.26.1–1.26.22.
- Turkmen A, Lin S. Are rare variants really independent? *Genet Epidemiol*. 2017 May; 41(4):363–71.
- Collins A, Morton NE. Mapping a disease locus by allelic association. *Proc Natl Acad Sci USA*. 1998 Feb;95(4):1741–5.

- 36 Shete S. A note on the optimal measure of allelic association. *Ann Hum Genet.* 2003 Mar; 67(Pt 2):189–91.
- 37 Mangin B, Garnier-Géré P, Cierco-Ayrolles C. The estimator of the optimal measure of allelic association: mean, variance and probability distribution when the sample size tends to infinity. *Stat Appl Genet Mol Biol.* 2008;7(1):20.
- 38 Weir BS. *Genetic Data Analysis II.* Sunderland: Sinauer; 1996.
- 39 Rosenberg NA, Blum MG. Sampling properties of homozygosity-based statistics for linkage disequilibrium. *Math Biosci.* 2007 Jul; 208(1):33–47.
- 40 Song YS, Song JS. Analytic computation of the expectation of the linkage disequilibrium coefficient r^2 . *Theor Popul Biol.* 2007 Feb;71(1): 49–60.
- 41 Alcalá N, Rosenberg NA. Mathematical constraints on F_{ST} : biallelic markers in arbitrarily many populations. *Genetics.* 2017 Jul;206(3): 1581–600.
- 42 Zapata C. The D' measure of overall gametic disequilibrium between pairs of multiallelic loci. *Evolution.* 2000 Oct;54(5):1809–12.
- 43 Payseur BA, Place M, Weber JL. Linkage disequilibrium between STRPs and SNPs across the human genome. *Am J Hum Genet.* 2008 May;82(5):1039–50.